



# Error estimate for finite volume approximate solutions of some oblique derivative boundary value problems

Abdallah Bradji, Thierry Gallouët

## ► To cite this version:

Abdallah Bradji, Thierry Gallouët. Error estimate for finite volume approximate solutions of some oblique derivative boundary value problems. *International Journal on Finite Volumes*, 2006, 3 (2), pp.1-35. hal-01114201

**HAL Id: hal-01114201**

**<https://hal.science/hal-01114201>**

Submitted on 9 Feb 2015

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Error estimate for finite volume approximate solutions of some oblique derivative boundary value problems

Abdallah BRADJI<sup>†</sup>

<sup>†</sup> *LATP, CMI, F-13453, Marseille*

bradji@cmi.univ-mrs.fr

Thierry GALLOUËT<sup>\*</sup>

<sup>\*</sup> *LATP, CMI, F-13453, Marseille*

gallouet@cmi.univ-mrs.fr

---

## Abstract

This paper is an improvement of [BG 05], concerning the Laplace equation with an oblique boundary condition. When the boundary condition involves a regular coefficient, we present a weak formulation of the problem and we prove some existence and uniqueness results of the weak solution. We develop a finite volume scheme and we prove the convergence of the finite volume solution to the weak solution, when the mesh size goes to zero. We also present some partial results for the interesting case of a discontinuous coefficient in the boundary condition. In particular, we give a finite volume scheme, taking in consideration the discontinuities of this coefficient. Finally, we obtain some error estimates (in a convenient norm) of order  $\sqrt{h}$  (where  $h$  is the mesh size), when the solution  $u$  is regular enough.

**Key words :** oblique derivative, smooth coefficient, piecewise constant coefficient, unstructured mesh, finite volume, error estimate.

---

## 1 Introduction

In this paper, we are interested with the finite volume approximation of the Laplace equation on an open bounded polygonal connected subset  $\Omega$  of  $\mathbb{R}^2$ , with an oblique boundary condition:

$$\begin{cases} -\Delta u(\mathbf{x}) = f(\mathbf{x}), & \mathbf{x} = (x, y) \in \Omega, \\ u_n(\mathbf{x}) + (\alpha u)_t(\mathbf{x}) = g(\mathbf{x}), & \mathbf{x} \in \partial\Omega, \end{cases} \quad (1)$$

with the notations  $v_n = \nabla v \cdot \mathbf{n}$  and  $v_t = \nabla v \cdot \mathbf{t}$ , where  $\mathbf{n} = (\mathbf{n}_x, \mathbf{n}_y)^t$  is the normal vector to the boundary  $\partial\Omega$ , outward to  $\Omega$ ,  $\mathbf{t} = (-\mathbf{n}_y, \mathbf{n}_x)^t$  (so that  $\mathbf{t}$  is a tangent vector to  $\partial\Omega$ ). The vectors  $\mathbf{n}$  and  $\mathbf{t}$  are defined everywhere on  $\partial\Omega$  except in a finite number of points. The functions  $f$ ,  $g$  and  $\alpha$  are given (in  $\Omega$  for  $f$  and on  $\partial\Omega$  for  $g$  and  $\alpha$ ).

In order to obtain an existence result for (1), a compatibility condition on  $f$  and  $g$  is needed. Indeed, assuming for simplicity that  $\alpha$  is a regular function and that  $u$  is a regular

solution to (1), this compatibility condition is easily obtained, taking a constant function as test function in (1), it reads:

$$\int_{\Omega} f(\mathbf{x}) d\mathbf{x} + \int_{\partial\Omega} g(\mathbf{x}) d\gamma(\mathbf{x}) = 0, \quad (2)$$

where  $\gamma$  stands for the one dimensional Lebesgue measure on  $\partial\Omega$ .

Similarly, in order to expect a uniqueness result for (1), we have to add an additional condition on  $u$ . For this additional condition, we will take:

$$\int_{\Omega} u(\mathbf{x}) d\mathbf{x} = 0. \quad (3)$$

A huge literature is devoted to the Laplace equation with various boundary conditions, namely Fourier, Neumann and Dirichlet boundary conditions, and to the discretization of such problems. A few papers are concerned by this quite unusual boundary condition, namely the oblique boundary condition given in (1). For instance, the Problem (1), in the case of  $\alpha$  is constant on each line of the boundary of  $\Omega$ , is considered in [G 85] and [M 74] (more precise, the operator considered in [M 74] is  $-\Delta u + u$  instead of  $-\Delta u$ ). They studied the regularity of the exact solutions by means of suitable a priori estimates.

The numerical study of oblique derivative boundary value problems is considered, for instance, in [M 02], where the author suggested the standard finite difference scheme of five points to approximate the solution of a nonlinear oblique derivative boundary value problem posed on a rectangular domain.

Problem (1) appears, for instance, in a method developed in [B 05] for improving the convergence order of numerical schemes for the classical Dirichlet problem (it can probably also appear in the modelization of some mechanical problems, but perhaps not directly under the form (1)). Instead of considering the boundary condition given in (1), it is also possible to consider  $u_n + \alpha u_t = g$ , which leads to the dual problem of (1). With this boundary condition the compatibility condition on  $f$  and  $g$  (in order to obtain an existence result) is not necessarily (2) and may depends on  $\alpha$ . On the contrary, the condition (3) is then quite natural since, for this problem, any constant function is solution with  $f = g = 0$ .

To discretize the problem (1), we introduce an unstructured mesh  $\mathcal{T}$  defined as in [EGH 00]. We consider three different cases for the problem (1) (with (3) and assuming (2)).

In the first case, we assume that  $\alpha$  is constant. We obtain a weak formulation and we present a finite volume scheme. We prove that the finite volume solution converges to the unique solution of the weak problem, when the mesh size goes to zero. If we assume that  $u \in \mathcal{C}^2(\overline{\Omega})$ , we prove that the error estimate is of order  $\sqrt{h}$ , where  $h$  is the size of  $\mathcal{T}$ .

In the second case, we assume that  $\alpha \in \mathcal{C}^1(\overline{\Omega})$ . Then, we also obtain an existence and uniqueness result of a weak solution when  $\alpha$  satisfies:

$$\min_{\partial\Omega} \alpha_t \geq -\delta,$$

where  $\delta$  is a positive real number only depending on  $\Omega$ . We present a finite volume solution and we prove its convergence to the weak solution of the problem. If we assume that the solution  $u$  belongs to  $\mathcal{C}^2(\overline{\Omega})$ , we give an error estimate of order  $\sqrt{h}$ .

Finally, we consider the case where  $\alpha$  is constant on each line of the boundary  $\partial\Omega$ . Then the oblique boundary condition of the problem (1) makes sense on each line of the boundary  $\partial\Omega$  (and can be written equivalently  $u_n + \alpha u_t = g$  or  $u_n + (\alpha u)_t = g$  on each line of the

boundary). Such a problem arises when  $u = v_x$  and  $v$  is a smooth function, satisfying the Laplace equation with a homogeneous boundary condition:

$$\begin{cases} -\Delta v(\mathbf{x}) = \bar{f}(\mathbf{x}), \mathbf{x} \in \Omega, \\ v(\mathbf{x}) = 0, \mathbf{x} \in \partial\Omega, \end{cases} \quad (4)$$

where  $f = \partial \bar{f} / \partial x$  and  $g$  (in (1)) is defined in terms of  $f$  and of the components of the tangential vector  $\mathbf{t}$  on the boundary  $\partial\Omega$ . Considering the equation satisfied by  $v_x$  is useful for improving the convergence order of numerical schemes for (4), see [A 06] and [B 05]. For this last case, we present a finite volume scheme, taking in consideration the discontinuities of  $\alpha$  on the corners of  $\Omega$ , and we prove an error estimate of order  $\sqrt{h}$ , assuming that  $u$  satisfies the Assumption 6.1 (note that if the solution  $v$  of the equation (4) belongs to  $\mathcal{C}^3(\bar{\Omega})$ , then  $u = v_x$  satisfies an equation of the form (1), in which  $\alpha$  is constant on each line of the boundary  $\partial\Omega$ , and  $u$  satisfies Assumption 6.1).

## 2 Preliminaries and functional spaces

Recall that the domain  $\Omega$  is an open bounded polygonal connected subset of  $\mathbb{R}^2$ . The boundary of  $\Omega$  is denoted by  $\partial\Omega$ . The norm in the usual Sobolev space  $H^1(\Omega)$  is defined by

$$\|w\|_{1,\Omega}^2 = \|w\|_{0,\Omega}^2 + \|\nabla w\|_{0,\Omega}^2,$$

where  $|\cdot|$  denotes the Euclidean norm in  $\mathbb{R}^2$  and  $\|\cdot\|_{0,\Omega}$  denotes the norm in the space  $L^2(\Omega)$  (a similar notation will be used for the norm in the space  $L^2(\partial\Omega)$ ). The space  $\mathcal{D}(\bar{\Omega})$  is the space of infinitely differentiable functions on  $\bar{\Omega}$  (that is to say the restrictions to  $\Omega$  of the infinitely differentiable functions defined on  $\mathbb{R}^2$ ) and  $\mathcal{D}(\Omega)$  as the space of infinitely differentiable functions, with compact support on  $\Omega$ . We recall that  $\mathcal{D}(\bar{\Omega})$  is a dense subspace of  $H^1(\Omega)$  and that the space  $H_0^1(\Omega)$  is the closure of  $\mathcal{D}(\Omega)$  in  $H^1(\Omega)$ . Thanks to the Lipchitz continuity of the boundary of  $\Omega$ , we also have  $H_0^1(\Omega) = \{v \in H^1(\Omega) : \tilde{\gamma}(v) = 0\}$ , where  $\tilde{\gamma}$  is the linear trace operator from  $H^1(\Omega)$  to  $L^2(\partial\Omega)$ .

Let  $\Gamma = \partial\Omega$ . We denote by  $H^{\frac{1}{2}}(\Gamma)$  the space of the “traces” on  $\Gamma$  of the elements of  $H^1(\Omega)$ , that is the Range of  $\tilde{\gamma}$ . The norm in  $H^{\frac{1}{2}}(\Gamma)$  is defined by:

$$\|w\|_{\frac{1}{2},\Gamma} = \inf_{\tilde{\gamma}(v)=w} \|v\|_{1,\Omega},$$

If  $u \in H^1(\Omega)$  and  $-\Delta u = f \in L^2(\Omega)$  (in the sense of  $\int_{\Omega} \nabla u \cdot \nabla v \, d\mathbf{x} = \int_{\Omega} f v \, d\mathbf{x}$ , for any  $v \in \mathcal{D}(\Omega)$  and then, by density of  $\mathcal{D}(\Omega)$  in  $H_0^1(\Omega)$ , for any  $v \in H_0^1(\Omega)$ ), we define the operator of normal derivative acting on  $u$  as the element  $u_n$  of  $H^{-\frac{1}{2}}(\Gamma)$  by:

$$\langle u_n, v \rangle_{H^{-\frac{1}{2}}(\Gamma), H^{\frac{1}{2}}(\Gamma)} = \int_{\Omega} \nabla u \cdot \nabla \tilde{v} \, d\mathbf{x} - \int_{\Omega} f \tilde{v} \, d\mathbf{x}, \forall v \in H^{\frac{1}{2}}(\Gamma). \quad (5)$$

If  $u \in H^1(\Omega)$ , we define the operator of tangential derivative acting on  $u$  as the element  $u_t$  of  $H^{-\frac{1}{2}}(\Gamma)$  by:

$$\langle u_t, v \rangle_{H^{-\frac{1}{2}}(\Gamma), H^{\frac{1}{2}}(\Gamma)} = \int_{\Omega} \tilde{v}_x u_y \, d\mathbf{x} - \int_{\Omega} u_x \tilde{v}_y \, d\mathbf{x}, \forall v \in H^{\frac{1}{2}}(\Gamma). \quad (6)$$

In (5) and (6),  $\tilde{v}$  is an element of  $H^1(\Omega)$  such that  $\tilde{\gamma}(\tilde{v}) = v$ . It is quite easy to see that these operators are well defined (indeed, thanks to a density argument and an integration by parts, the right-hand-sides of (5) and (6) vanish for any  $\tilde{v} \in H_0^1(\Omega)$ , and then they do not depend on the choice of  $\tilde{v}$  provided that  $\tilde{\gamma}(\tilde{v}) = v$ ). Using integration by parts, it is also

easy to see that  $u_n$  and  $u_t$  correspond to the classical derivatives of  $u$ ,  $\nabla u \cdot \mathbf{n}$  and  $\nabla u \cdot \mathbf{t}$ , when  $u$  is a regular function.

The definitions of the normal and tangential derivative operators given by (5) and (6) enable us to give a sense to the boundary condition in (1) when  $f \in L^2(\Omega)$ ,  $g \in L^2(\partial\Omega)$  (or  $g \in H^{-\frac{1}{2}}(\partial\Omega)$ ) and when  $\alpha$  is a constant function or, more generally, when  $\alpha$  is a smooth function (say  $\alpha \in C^1(\partial\Omega)$ , that is to say the restriction to  $\partial\Omega$  of a  $C^1$  function defined on  $\mathbb{R}^2$ ). Indeed, these definitions give a weak formulation for the problem (1) when  $\alpha$  is constant (see Theorem 4.2) and when  $\alpha$  is a smooth function (see Theorem 5.2).

### 3 Finite volume meshes

We first describe the assumptions which are needed on the mesh.

**DEFINITION 3.1** (Admissible meshes, cf. Eymard *et al.* [EGH 00]) An admissible finite volume mesh of  $\Omega$ , denoted by  $\mathcal{T}$ , is a finite family of open polygonal convex disjoint subsets of  $\Omega$  (the “control volumes”), with positive measures. To this family is associated a family of disjoint subsets of  $\overline{\Omega}$  contained in hyperplanes of  $\mathbb{R}^2$ , denoted by  $\mathcal{E}$  (these are the edges of the control volumes) and a family of points of  $\Omega$ ,  $\mathcal{P} = \{\mathbf{x}_K, K \in \mathcal{T}\}$ , satisfying the following properties (as in Definition 9.1 of [EGH 00]):

- $\overline{\Omega} = \cup_{K \in \mathcal{T}} \overline{K}$ . For any  $K \in \mathcal{T}$ , let  $\partial K = \overline{K} \setminus K$  be the boundary of  $K$ . For all  $K \in \mathcal{T}$ ,  $m(K)$  is the two-dimensional Lebesgue measure of  $K$  (it is the area of  $K$ ).
- For all  $\sigma \in \mathcal{E}$ , there exists a hyperplane  $E$  of  $\mathbb{R}^2$  and  $K \in \mathcal{T}$  with  $\overline{\sigma} = \partial K \cap E$  and  $\sigma$  is a non empty open subset of  $E$ . We then denote by  $m(\sigma)$  the one dimensional measure of  $\sigma$  and one assumes  $m(\sigma) > 0$ . We assume that, for all  $K \in \mathcal{T}$ , there exists a subset  $\mathcal{E}_K$  of  $\mathcal{E}$  such that  $\partial K = \cup_{\sigma \in \mathcal{E}_K} \overline{\sigma}$ . It then results from the previous hypotheses that, for all  $\sigma \in \mathcal{E}$ , either  $\sigma \subset \partial\Omega$  or there exists  $(K, L) \in \mathcal{T}^2$  with  $K \neq L$  such that  $\overline{K} \cap \overline{L} = \overline{\sigma}$ ; we denote in the latter case  $\sigma = K|L$ .
- For all  $K \in \mathcal{T}$ ,  $\mathbf{x}_K \in K$ . Furthermore, for all  $\sigma \in \mathcal{E}$  such that there exists  $(K, L) \in \mathcal{T}^2$  with  $\sigma = K|L$ , it is assumed that the straight line  $(\mathbf{x}_K, \mathbf{x}_L)$  going through  $\mathbf{x}_K$  and  $\mathbf{x}_L$  is orthogonal to  $K|L$ . For  $K \in \mathcal{T}$  and  $\sigma \in \mathcal{E}_K$ , let  $\mathcal{D}_{K,\sigma}$  be the straight line going through  $\mathbf{x}_K$  and orthogonal to  $\sigma$ . We assume that  $\mathcal{D}_{K,\sigma} \cap \sigma \neq \emptyset$  and we set  $\{\mathbf{y}_\sigma\} = \mathcal{D}_{K,\sigma} \cap \sigma$ .

If  $\mathcal{T}$  is an admissible mesh, we will also use the following notations:

- The mesh size is defined by  $\text{size}(\mathcal{T}) = \sup\{\text{diam}(K), K \in \mathcal{T}\}$ ,
- the set of interior (resp. boundary) edges is denoted by  $\mathcal{E}_{\text{int}}$  (resp.  $\mathcal{E}_{\text{ext}}$ ),  $\mathcal{E}_{\text{int}} = \{\sigma \in \mathcal{E} : \sigma \not\subset \partial\Omega\}$  (resp.  $\mathcal{E}_{\text{ext}} = \{\sigma \in \mathcal{E} : \sigma \subset \partial\Omega\}$ ),
- the set of neighbours of  $K$  is denoted by  $\mathcal{N}(K)$ ,  $\mathcal{N}(K) = \{L \in \mathcal{T} : \exists \sigma \in \mathcal{E}_K, \overline{\sigma} = \overline{K} \cap \overline{L}\}$ ,
- if  $\sigma = K|L$ , we denote by  $d_\sigma$  or  $d_{K|L}$  the Euclidean distance between  $\mathbf{x}_K$  and  $\mathbf{x}_L$  (which is positive) and  $d_{K,\sigma}$  the distance from  $\mathbf{x}_K$  to  $\sigma$ ,
- if  $\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}$ , let  $d_\sigma$  denote the Euclidean distance between  $\mathbf{x}_K$  and  $\mathbf{y}_\sigma$  (then,  $d_\sigma = d_{K,\sigma}$ ),
- for any  $\sigma \in \mathcal{E}$ , the “transmissibility” through  $\sigma$  is defined by  $\tau_\sigma = \frac{m(\sigma)}{d_\sigma}$  (note that  $d_\sigma > 0$ ).

To discretize the oblique boundary condition (the second equation of the problem (1)), we need the following definition:

**DEFINITION 3.2** Let  $\sigma \in \mathcal{E}_{\text{ext}}$  and  $\mathbf{n}$  be the normal vector to  $\sigma$ , outward to  $\Omega$ . Recall that  $\mathbf{t} = (-\mathbf{n}_y, \mathbf{n}_x)^t$  where  $\mathbf{n} = (\mathbf{n}_x, \mathbf{n}_y)^t$ . Then  $\sigma = (a, b) = \{sa + (1-s)b, s \in (0, 1)\}$  where  $a, b$  are chosen such that  $|b-a|\mathbf{t} = b-a$ . We denote by  $\sigma^-$  (resp.  $\sigma^+$ ) the element of  $\mathcal{E}_{\text{ext}}$  such that  $a$  is in the closure of  $\sigma^-$  (resp.  $b$  is in the closure of  $\sigma^+$ ) and  $\sigma^- \neq \sigma$  (resp.  $\sigma^+ \neq \sigma$ ). We also set  $\sigma_e = b$  and  $\sigma_b = a$  (so that  $|\sigma_e - \sigma_b|\mathbf{t} = \sigma_e - \sigma_b$ ).

For the discrete solution, we use the following space:

**DEFINITION 3.3** (The Finite volume space) For an admissible mesh  $\mathcal{T}$ , the space  $\mathcal{X}(\mathcal{T})$  is defined by  $\mathcal{X}(\mathcal{T}) = \mathcal{Y}(\mathcal{T}) \times \mathcal{Z}(\mathcal{T}) \subset L^2(\Omega) \times L^2(\partial\Omega)$  where  $\mathcal{Y}(\mathcal{T})$  is the set of functions from  $\Omega$  to  $\mathbb{R}$  which are constant over each control volume  $K \in \mathcal{T}$  and  $\mathcal{Z}(\mathcal{T})$  be the set of functions which are constant on each  $\sigma \in \mathcal{E}_{\text{ext}}$ . Thus an element of  $\mathcal{X}(\mathcal{T})$  is of the form  $(u_{\mathcal{T}}, v_{\mathcal{T}}) \in L^2(\Omega) \times L^2(\partial\Omega)$  where  $u_{\mathcal{T}}$  (resp.  $v_{\mathcal{T}}$ ) is constant over each control volume  $K \in \mathcal{T}$  (resp. constant over each boundary edge  $\sigma \in \mathcal{E}_{\text{ext}}$ ).

To analyze the convergence of the finite volume schemes, we use the following semi-norm on  $\mathcal{X}(\mathcal{T})$ :

**DEFINITION 3.4** (Discrete semi-norm on  $\mathcal{X}(\mathcal{T})$ ) Let  $(u_{\mathcal{T}}, v_{\mathcal{T}}) \in \mathcal{X}(\mathcal{T})$  ( $\mathcal{X}(\mathcal{T})$  given by Definition 3.3), one defines a discrete semi-norm by:

$$|(u_{\mathcal{T}}, v_{\mathcal{T}})|_{1, \mathcal{X}(\mathcal{T})}^2 = |u_{\mathcal{T}}|_{1, \mathcal{T}}^2 + \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} \frac{m(\sigma)}{d_{\sigma}} (u_K - u_{\sigma})^2, \quad (7)$$

where

$$|u|_{1, \mathcal{T}}^2 = \sum_{\sigma \in \mathcal{E}_{\text{int}}} \frac{m(\sigma)}{d_{\sigma}} (D_{\sigma} u)^2, \quad (8)$$

$D_{\sigma} u = |u_L - u_K|$ , if  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $\sigma = K|L$ , and  $u_K$  (resp.  $u_{\sigma}$ ) denotes the value taken by  $u_{\mathcal{T}}$  (resp.  $v_{\mathcal{T}}$ ) on the control volume  $K$  (resp. on the boundary edge  $\sigma$ ).

*Remark 1* With the notations of Definition 3.4, let  $(u_{\mathcal{T}}, v_{\mathcal{T}}) \in \mathcal{X}(\mathcal{T})$  such that:

$$|(u_{\mathcal{T}}, v_{\mathcal{T}})|_{1, \mathcal{X}(\mathcal{T})} = 0.$$

Then, for all  $\sigma \in \mathcal{E}_{\text{int}}$ , one has  $u_K = u_L$  where  $\sigma = K|L$  (since  $m(\sigma)/d_{\sigma} > 0$ ). Since  $\Omega$  is connected, we then deduce that  $u_{\mathcal{T}}$  is a constant function. For all  $\sigma \in \mathcal{E}_{\text{ext}}$ , one also has  $u_K = u_{\sigma}$  where  $\sigma \in \mathcal{E}_K$  (since  $m(\sigma)/d_{\sigma} > 0$ ). Then, there exists  $C \in \mathbb{R}$  such that  $u_K = C$  for all  $K \in \mathcal{T}$  and  $u_{\sigma} = C$  for all  $\sigma \in \mathcal{E}_{\text{ext}}$ .

To define a finite volume scheme for (1), we consider an admissible mesh  $\mathcal{T}$  in the sense of Definition 3.1 and we use the following quantities:

$$f_K = \frac{1}{m(K)} \int_K f(\mathbf{x}) d\mathbf{x} \text{ and } g_{\sigma} = \frac{1}{m(\sigma)} \int_{\sigma} g(\mathbf{x}) d\gamma(\mathbf{x}). \quad (9)$$

In Problem (1), there are an equation on the domain  $\Omega$  and an equation on the boundary  $\partial\Omega$ . To get a finite volume scheme, we integrate the first equation on the control volumes and the second one on the boundary edges. The first integration can be done using the usual techniques of [EGH 00]. For the second integration, we use the following useful property. Let  $a$  and  $b$  be two points in  $\mathbb{R}^2$  and  $(a, b) = \{sa + (1-s)b, s \in (0, 1)\}$ . Let  $f \in \mathcal{C}^1(\mathbb{R}^2)$  and  $\mathbf{t} = \frac{b-a}{|b-a|}$ . Let  $f_t = \nabla f \cdot \mathbf{t}$  (it is the tangential derivative of  $f$  on  $(a, b)$ ). Then:

$$\int_{(a,b)} f_t(\mathbf{x}) d\gamma(\mathbf{x}) = f(b) - f(a). \quad (10)$$

## 4 The case “ $\alpha$ constant”

In this Section, we are interested by Problem (1) when the function  $\alpha$  is constant. Then, Problem (1) can be rewritten, with  $\alpha \in \mathbb{R}$ , as:

$$\begin{cases} -\Delta u(\mathbf{x}) = f(\mathbf{x}), & \mathbf{x} \in \Omega, \\ u_n(\mathbf{x}) + \alpha u_t(\mathbf{x}) = g(\mathbf{x}), & \mathbf{x} \in \partial\Omega, \end{cases} \quad (11)$$

*Remark 2* The problem (11) can be rewritten as a Neumann Problem (with, if  $\alpha \neq 0$ , a non symmetric operator):

$$\begin{cases} -\operatorname{div}(\mathcal{A} \operatorname{grad} u(\mathbf{x})) = f(\mathbf{x}), & \mathbf{x} \in \Omega, \\ (\mathcal{A} \operatorname{grad} u(\mathbf{x})) \cdot \mathbf{n}(\mathbf{x}) = g(\mathbf{x}), & \mathbf{x} \in \Gamma, \end{cases} \quad (12)$$

where the positive definite matrix  $\mathcal{A}$  is given by:

$$\mathcal{A} = \begin{pmatrix} 1 & \alpha \\ -\alpha & 1 \end{pmatrix}$$

But the finite volume scheme we shall present will be derived directly from the equation (11) and not from (12) (see Section 4.1 below and Sections 10 and 11 in [EGH 00]).

To get the existence of a solution for problem (11), we assume a compatibility condition on  $f$  and  $g$  (see (2)):

**ASSUMPTION 4.1**  $(g, f) \in L^2(\partial\Omega) \times L^2(\Omega)$  and  $\int_{\Omega} f(\mathbf{x}) d\mathbf{x} + \int_{\partial\Omega} g(\mathbf{x}) d\gamma(\mathbf{x}) = 0$ , where  $\gamma$  is one dimensional Lebesgue measure on  $\partial\Omega$ .

Under this assumption, the following theorem yields existence and uniqueness of a weak solution to Problem (11) (or (1)) with Condition (3). The proof of this theorem is an easy consequence of the Lax-Milgram lemma. This proof is given in the more general case of a smooth function  $\alpha$  in Section 5 (see Theorem 5.2)

**THEOREM 4.2** Let  $\alpha \in \mathbb{R}$ . Under the Assumption 4.1, there exists a unique solution of (13)-(14):

$$u \in H^1(\Omega), \int_{\Omega} u(\mathbf{x}) d\mathbf{x} = 0, \quad (13)$$

$$\int_{\Omega} \nabla u \cdot \nabla v d\mathbf{x} + \alpha \int_{\Omega} (v_x u_y - u_x v_y) d\mathbf{x} = \int_{\Omega} f v d\mathbf{x} + \int_{\partial\Omega} g \tilde{\gamma}(v) d\gamma(\mathbf{x}), \forall v \in H^1(\Omega). \quad (14)$$

Let  $\alpha \in \mathbb{R}$ ,  $f \in L^2(\Omega)$  and  $g \in L^2(\partial\Omega)$ . Thanks to the definitions of Section 2, let  $u$  be a function satisfying (13), then,  $u$  is a solution of (14) if and only if  $u$  satisfies:

- $-\Delta u = f$  in  $\mathcal{D}'(\Omega)$  (the usual dual space of  $\mathcal{D}(\Omega)$ ),
- $u_n + \alpha u_t = g$  in  $H^{-\frac{1}{2}}(\partial\Omega)$ , with, for  $w \in H^{\frac{1}{2}}(\partial\Omega) (\subset L^2(\partial\Omega))$ ,

$$\langle g, w \rangle_{H^{-\frac{1}{2}}(\partial\Omega), H^{\frac{1}{2}}(\partial\Omega)} = \int_{\partial\Omega} g(\mathbf{x}) w(\mathbf{x}) d\gamma(\mathbf{x}).$$

#### 4.1 The finite volume scheme for (11)

*Remark 3* For the sake of simplicity, we assume that  $\alpha > 0$  in (11). For  $\alpha = 0$ , we get the classical Neumann problem which is treated in [EGH 00]. The case  $\alpha < 0$  can be treated in a similar way to the case  $\alpha > 0$ .

In order to obtain the numerical scheme we will discretize (11) (or (1)) instead of, for instance, (12). This choice allows the use of an “admissible mesh” for the Laplace equation (instead of an “admissible mesh” for the matrix  $\mathcal{A}$ ) and allows the resolution of the discrete problem using the matrix corresponding to the discretization of the Laplace equation with a homogeneous Dirichlet boundary condition (see, for instance, [B 05]).

Let  $(u_K)_{K \in \mathcal{T}}$  and  $(u_\sigma)_{\sigma \in \mathcal{E}_{\text{ext}}}$  denote the discrete unknowns. The numerical scheme is defined by the following set of equations:

$$\sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma} = m(K)f_K, \forall K \in \mathcal{T}, \quad (15)$$

where

$$F_{K,\sigma} = -\tau_{K|L}(u_L - u_K), \forall \sigma \in \mathcal{E}_{\text{int}}, \text{ if } \sigma = K|L, \quad (16)$$

$$F_{K,\sigma} = -\tau_\sigma(u_\sigma - u_K), \forall \sigma \in \mathcal{E}_{\text{ext}} \text{ such that } \sigma \in \mathcal{E}_K, \quad (17)$$

and

$$\tau_\sigma(u_\sigma - u_K) = -\alpha(u_\sigma - u_{\sigma^-}) + m(\sigma)g_\sigma, \forall \sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}, \quad (18)$$

where  $\sigma^-$  is defined in Definition 3.2.

The condition (13) can be discretized by

$$\sum_{K \in \mathcal{T}} m(K)u_K = 0. \quad (19)$$

*Remark 4* The unknowns  $\{u_K, K \in \mathcal{T}\}$  and  $\{u_\sigma, \sigma \in \mathcal{E}_{\text{ext}}\}$  of the finite volume scheme are expected to approximate  $u$  on the control volumes  $\{K\}_{K \in \mathcal{T}}$  through  $\{u_K\}_{K \in \mathcal{T}}$ , and are expected to approximate  $u$  on  $\{\sigma\}_{\sigma \in \mathcal{E}_{\text{ext}}}$  (see Definition 3.1) through  $\{u_\sigma\}_{\sigma \in \mathcal{E}_{\text{ext}}}$ .

With the discrete unknowns, namely  $\{u_K, K \in \mathcal{T}\}$  and  $\{u_\sigma, \sigma \in \mathcal{E}_{\text{ext}}\}$ , it is possible to define an element of  $\mathcal{X}(\mathcal{T})$ . If the discrete unknowns satisfy (15)-(19), we will say that this element of  $\mathcal{X}(\mathcal{T})$  is a solution of (15)-(19):

**DEFINITION 4.3** An element  $(u_{\mathcal{T}}, v_{\mathcal{T}}) \in \mathcal{X}(\mathcal{T})$  (see Definition 3.3) is a solution of (15)-(19) if  $u_{\mathcal{T}}(\mathbf{x}) = u_K$  for  $\mathbf{x} \in K$ , for all  $K \in \mathcal{T}$ , and  $v_{\mathcal{T}}(\mathbf{x}) = u_\sigma$  for  $\mathbf{x} \in \sigma$ , for all  $\sigma \in \mathcal{E}_{\text{ext}}$ , where  $(u_K)_{K \in \mathcal{T}}$  and  $(u_\sigma)_{\sigma \in \mathcal{E}_{\text{ext}}}$  satisfy (15)-(19).

#### 4.2 Existence and uniqueness of the discrete solution

We use the techniques of the Proof of Lemma 10.1 in [EGH 00] and the following equality to prove the existence and uniqueness of the solution of (15)-(19):

$$\sum_{\sigma \in \mathcal{E}_{\text{ext}}} (u_\sigma - u_{\sigma^-})u_\sigma = \frac{1}{2} \sum_{\sigma \in \mathcal{E}_{\text{ext}}} (u_\sigma - u_{\sigma^-})^2. \quad (20)$$

**THEOREM 4.4** Let  $\alpha \in \mathbb{R}$ . Assume that  $\alpha > 0$  (for the other cases, see Remark 3). Let  $\mathcal{T}$  be an admissible mesh in the sense of Definition 3.1 and  $\{(f_K, g_\sigma), (K, \sigma) \in \mathcal{T} \times \mathcal{E}_{\text{ext}}\}$  defined by (9). Then, under the Assumption 4.1, the system (15)-(19) has a unique solution.



Let  $(u_{\mathcal{T}}, v_{\mathcal{T}}) \in \mathcal{X}(\mathcal{T})$  be the solution of (15)-(19) in the sense of Definition 4.3. Then, there exists  $C_1 \in \mathbb{R}^+$ , only depending on  $\Omega$ , such that

$$\left( |(u_{\mathcal{T}}, v_{\mathcal{T}})|_{1, \mathcal{X}(\mathcal{T})}^2 + \frac{\alpha}{2} |v_{\mathcal{T}}|_{\mathcal{Z}(\mathcal{T})}^2 \right)^{\frac{1}{2}} \leq C_1 (\|f\|_{0, \Omega} + \|g\|_{0, \partial\Omega}), \quad (21)$$

where the semi-norm  $|(u_{\mathcal{T}}, v_{\mathcal{T}})|_{1, \mathcal{X}(\mathcal{T})}$  is defined in Definition 3.4 and  $|v_{\mathcal{T}}|_{\mathcal{Z}(\mathcal{T})}$  is the semi-norm defined by:

$$|v_{\mathcal{T}}|_{\mathcal{Z}(\mathcal{T})}^2 = \sum_{\sigma \in \mathcal{E}_{\text{ext}}} (u_{\sigma} - u_{\sigma-})^2. \quad (22)$$

### Proof

**Step 1.** Existence and uniqueness of the solution of (15)-(19).

Let  $\{u_K, K \in \mathcal{T}\}$  and  $\{u_{\sigma}, \sigma \in \mathcal{E}_{\text{ext}}\}$  be a solution of (15)-(18). Multiplying both sides of Equation (15) by  $u_K$ , summing over  $K, K \in \mathcal{T}$ , using (16)-(18) and (20), we get:

$$|(u_{\mathcal{T}}, v_{\mathcal{T}})|_{1, \mathcal{X}(\mathcal{T})}^2 + \frac{\alpha}{2} |v_{\mathcal{T}}|_{\mathcal{Z}(\mathcal{T})}^2 = \mathbb{T}_1^{\mathcal{T}} + \mathbb{T}_2^{\mathcal{T}}, \quad (23)$$

where

$$\mathbb{T}_1^{\mathcal{T}} = \sum_{\sigma \in \mathcal{E}_{\text{ext}}} m(\sigma) g_{\sigma} u_{\sigma} \text{ and } \mathbb{T}_2^{\mathcal{T}} = \sum_{K \in \mathcal{T}} m(K) f_K u_K. \quad (24)$$

Let  $M_1$  be the number of elements of  $\mathcal{T}$ ,  $M_2$  the number of elements of  $\mathcal{E}_{\text{ext}}$  and  $M = M_1 + M_2$ . The system (15)-(18) can be viewed as a system of  $M$  unknowns (which are  $\{u_K, K \in \mathcal{T}\}$  and  $\{u_{\sigma}, \sigma \in \mathcal{E}_{\text{ext}}\}$ ) with  $M$  equations. After the choice of an order for the unknowns and the equations, it can be written as  $Aw = b$ , where  $A$  is  $M \times M$  matrix,  $w \in \mathbb{R}^M$  is the unknown vector and  $b \in \mathbb{R}^M$  is given by the data (namely  $f$  and  $g$ ). Equality (23) proves that if  $b = 0$  (that is  $f_K = 0$  for all  $K \in \mathcal{T}$  and  $g_{\sigma} = 0$  for all  $\sigma \in \mathcal{E}_{\text{ext}}$ ) then  $|(u_{\mathcal{T}}, v_{\mathcal{T}})|_{1, \mathcal{X}(\mathcal{T})} = 0$ . Following Remark 1, one deduces that there exists  $C \in \mathbb{R}$  such that  $u_K = C$  for all  $K \in \mathcal{T}$  and  $u_{\sigma} = C$  for all  $\sigma \in \mathcal{E}_{\text{ext}}$ . This proves that the dimension of the null space of  $A$  is 1. Therefore, the dimension of the range of  $A$  is  $M - 1$ . Since  $\sum_{K \in \mathcal{T}} m(K) f_K + \sum_{\sigma \in \mathcal{E}_{\text{ext}}} m(\sigma) g_{\sigma} = 0$  is a necessary condition for (15)-(18) to have a solution, it is also a sufficient condition. Furthermore, under this condition on  $f$  and  $g$  (which is given by Assumption 4.1), since the null space of  $A$  is reduced to the set of constant vectors, the system (15)-(19) has a unique solution.

**Step 2.** Proof of Estimate (21).

Under Assumption 4.1, let  $(u_{\mathcal{T}}, v_{\mathcal{T}}) \in \mathcal{X}(\mathcal{T})$  be the solution of (15)-(19) in the sense of Definition 4.3. Using (24), the Cauchy Schwarz inequality and inequalities (10.10) (the discrete mean Poincaré inequality) and (10.25) (the discrete trace inequality) of [EGH 00] (combined with equation (19)), we get (using also  $d_{\sigma} \leq \text{diam}(\Omega)$ ):

$$|\mathbb{T}_1^{\mathcal{T}}| \leq C_2 \left( \left( \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} \frac{m(\sigma)}{d_{\sigma}} (u_K - u_{\sigma})^2 \right)^{\frac{1}{2}} + |u_{\mathcal{T}}|_{1, \mathcal{T}} \right) \|g\|_{0, \partial\Omega}, \quad (25)$$

and

$$|\mathbb{T}_2^{\mathcal{T}}| \leq C_3 \|f\|_{0, \Omega} |u_{\mathcal{T}}|_{1, \Omega}, \quad (26)$$

where  $C_2$  and  $C_3$  are only depending on  $\Omega$ .

Combining (23), (25) and (26) yields (21). ■

### 4.3 The convergence of $(u_{\mathcal{T}}, v_{\mathcal{T}})$

In this section, we prove the convergence of the solution  $(u_{\mathcal{T}}, v_{\mathcal{T}})$  of (15)-(19) when the size of the mesh goes to 0, assuming Assumption 4.1 and the additional assumption (34), for some  $\zeta_1 > 0$ . The proof is mainly based on the ideas developed in [EGH 00] for the Neumann problem. In all this section, we assume that  $\alpha > 0$  (for the other cases, see Remark 3),  $\mathcal{T}$  is an admissible mesh in the sense of Definition 3.1 and that  $\{(f_K, g_\sigma), (K, \sigma) \in \mathcal{T} \times \mathcal{E}_{\text{ext}}\}$  is defined by (9). Assuming Assumption 4.1, the system (15)-(19) has a unique solution (thanks to Theorem 4.4). We first prove the following Lemma:

**LEMMA 4.5** Let  $(u_{\mathcal{T}}, v_{\mathcal{T}}) \in \mathcal{X}(\mathcal{T})$  be the solution of (15)-(19) in the sense of Definition 4.3. Then, there exists  $C_4$ , only depending on  $(\Omega, f, g)$ , such that

$$\|v_{\mathcal{T}}\|_{0,\partial\Omega} \leq C_4. \quad (27)$$

**Proof** Thanks to inequalities (10.25) and (10.10) of [EGH 00], (19) and (21), there exists a constant  $C_5$ , only depending on  $\Omega$ , such that

$$\begin{aligned} \sum_{\sigma \in \mathcal{E}_{\text{ext}}} m(\sigma) u_\sigma^2 &\leq 2 \left( \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} m(\sigma) u_K^2 + \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} m(\sigma) (u_\sigma - u_K)^2 \right) \\ &\leq C_5 \left( |u_{\mathcal{T}}|_{1,\mathcal{T}}^2 + |(u_{\mathcal{T}}, v_{\mathcal{T}})|_{1,\mathcal{X}(\mathcal{T})}^2 \right) \leq C_4, \end{aligned}$$

where  $C_4 = 2 C_5 C_1^2 (\|f\|_{0,\Omega} + \|g\|_{0,\partial\Omega})^2$ . ■

Since the set  $Y$  of the approximations  $u_{\mathcal{T}}$  is bounded in  $L^2(\Omega)$  (thanks to discrete mean Poincaré inequality and inequality (21)), we are able now to justify that  $u_{\mathcal{T}}$  converges to some  $u$  as  $\text{size}(\mathcal{T})$  goes to 0.

Uniform boundedness (21) and compactness result of [EGH 00] in case of Neumann problem yield that the set  $Y$  is relatively compact in  $L^2(\Omega)$ . In addition, if a sequence  $u_{\mathcal{T}_n}$  converges to a function  $u$  in  $L^2$ -norm as  $\text{size}(\mathcal{T}_n)$  goes to 0, then  $u \in H^1(\Omega)$ . Furthermore, Lemma 4.5 implies that  $v_{\mathcal{T}_n}$  converges weakly to some  $v \in L^2(\partial\Omega)$ , up to a subsequence. We start by proving that:

$$\begin{aligned} - \int_{\Omega} u(\mathbf{x}) \Delta \varphi(\mathbf{x}) d\mathbf{x} + \int_{\partial\Omega} \varphi_n(\mathbf{x}) v(\mathbf{x}) d\gamma(\mathbf{x}) &= \int_{\Omega} f(\mathbf{x}) \varphi(\mathbf{x}) d\mathbf{x} \\ + \int_{\partial\Omega} g(\mathbf{x}) \varphi(\mathbf{x}) d\gamma(\mathbf{x}) + \alpha \int_{\partial\Omega} \varphi_t(\mathbf{x}) v(\mathbf{x}) d\gamma(\mathbf{x}), \quad \forall \varphi \in \mathcal{C}^2(\overline{\Omega}). \end{aligned} \quad (28)$$

To simplify the notations, we set  $u_{\mathcal{T}_n} = u_{\mathcal{T}}$  and  $v_{\mathcal{T}_n} = v_{\mathcal{T}}$ . Let  $\varphi \in \mathcal{C}^2(\overline{\Omega})$  and consider the function  $\varphi_{\mathcal{T}} = (\varphi_{\mathcal{T}}^{(1)}, \varphi_{\mathcal{T}}^{(2)}) \in \mathcal{X}(\mathcal{T})$  (see Definition 3.3) defined by  $\varphi_{\mathcal{T}}^{(1)}(\mathbf{x}) = \varphi_K = \varphi(\mathbf{x}_K)$ , for  $\mathbf{x} \in K$  and for any control volume  $K$ , and  $\varphi_{\mathcal{T}}^{(2)}(\mathbf{x}) = \varphi_\sigma = \varphi(\mathbf{y}_\sigma)$  for  $\mathbf{x} \in \sigma$ , for any  $\sigma \in \mathcal{E}_{\text{ext}}$  (see Definition 3.1). Multiplying both sides of equation (15) by  $\varphi_K$ , summing over  $K \in \mathcal{T}$  and reordering the terms yields

$$- \sum_{K \in \mathcal{T}} u_K \sum_{L \in \mathcal{N}(K)} \tau_{L|K} (\varphi_L - \varphi_K) = \int_{\Omega} f(\mathbf{x}) \varphi_{\mathcal{T}}^{(1)}(\mathbf{x}) d\mathbf{x} + \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} m(\sigma) \frac{u_\sigma - u_K}{d_\sigma} \varphi_K. \quad (29)$$

Using the consistency of the flux and the fact that  $\varphi$  is smooth one has (cf. [EGH 00], page 813):

$$\sum_{L \in \mathcal{N}(K)} \tau_{L|K} (\varphi_L - \varphi_K) = \int_K \Delta \varphi(\mathbf{x}) d\mathbf{x} - \int_{\partial\Omega \cap \partial K} \varphi_n(\mathbf{x}) d\gamma(\mathbf{x}) + \sum_{L \in \mathcal{N}(K)} R_{K,L}(\varphi), \quad (30)$$

with  $R_{K,L} = -R_{L,K}$ , for all  $L \in \mathcal{N}(K)$  and  $K \in \mathcal{T}$ , and  $|R_{K,L}| \leq C_6 \mathfrak{m}(K|L) \text{size}(\mathcal{T})$ , where  $C_6$  only depends on  $\varphi$ . This with (29) and (18) implies:

$$\begin{aligned} - \int_{\Omega} u_{\mathcal{T}}(\mathbf{x}) \Delta \varphi(\mathbf{x}) d\mathbf{x} &+ \int_{\partial\Omega} \varphi_n(\mathbf{x}) v_{\mathcal{T}}(\mathbf{x}) d\gamma(\mathbf{x}) + \bar{r} = \int_{\Omega} f(\mathbf{x}) \varphi_{\mathcal{T}}^{(1)}(\mathbf{x}) d\mathbf{x} \\ &+ \int_{\partial\Omega} g(\mathbf{x}) \varphi_{\mathcal{T}}^{(2)}(\mathbf{x}) d\gamma(\mathbf{x}) - \alpha \sum_{\sigma \in \mathcal{E}_{\text{ext}}} (u_{\sigma} - u_{\sigma-}) \varphi_{\sigma} \\ &- \alpha \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} (u_{\sigma} - u_{\sigma-}) (\varphi_K - \varphi_{\sigma}) \\ &+ \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} \mathfrak{m}(\sigma) (\varphi_K - \varphi_{\sigma}) g_{\sigma}, \end{aligned} \quad (31)$$

where  $\bar{r} = r(\varphi, \mathcal{T}) + s(\varphi, \mathcal{T})$ , with  $r(\varphi, \mathcal{T}) = -\sum_{K \in \mathcal{T}} u_K \sum_{L \in \mathcal{N}(K)} R_{K,L}$  and  $s$  is given by

$$s(\varphi, \mathcal{T}) = \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} (u_K - u_{\sigma}) \int_{\sigma} \varphi_n(\mathbf{x}) d\gamma(\mathbf{x}) \quad (32)$$

In order to pass to the limit in (31) as  $\text{size}(\mathcal{T})$  goes to 0, we need some estimates. Reordering the sum in  $r(\varphi, \mathcal{T})$  and using (21), we obtain

$$|r(\varphi, \mathcal{T})| \leq C_7 \text{size}(\mathcal{T}), \quad (33)$$

where  $C_7$  is a real positive number only depending on  $\Omega, f, g$  and  $\varphi$ .

On the other hand, if the mesh  $\mathcal{T}$  satisfies, for some  $\zeta_1 > 0$ :

$$d_{\sigma} \leq \zeta_1 \mathfrak{m}(\sigma), \forall \sigma \in \mathcal{E}_{\text{ext}}, \quad (34)$$

then, by using Cauchy Schwarz inequality, the uniform boundedness (21) and the regularity of  $\varphi$ , there exist  $C_8$  only depending on  $(\Omega, \alpha, \zeta_1, f, g)$ ,  $C_9$  only depending on  $(\Omega, f, g)$ ,  $C_{10}$  only depending on  $\Omega$  and  $C_{11}$  only depending on  $(\Omega, \alpha, f, g)$  such that

$$\left| \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} (u_{\sigma} - u_{\sigma-}) (\varphi_K - \varphi_{\sigma}) \right| \leq C_8 \sqrt{\text{size}(\mathcal{T})} |\varphi|_{1,\infty,\bar{\Omega}}, \quad (35)$$

$$|s(\varphi, \mathcal{T})| \leq C_9 \sqrt{\text{size}(\mathcal{T})} \|\varphi\|_{2,\Omega}, \quad (36)$$

$$\left| \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} \mathfrak{m}(\sigma) (\varphi_K - \varphi_{\sigma}) g_{\sigma} \right| \leq C_{10} \text{size}(\mathcal{T}) \|g\|_{0,\partial\Omega} \|\varphi\|_{1,\infty,\bar{\Omega}}, \quad (37)$$

$$\left| \sum_{\sigma \in \mathcal{E}_{\text{ext}}} (\varphi(\mathbf{y}_{\sigma}) - \varphi(\sigma_b)) (u_{\sigma} - u_{\sigma-}) \right| \leq C_{11} \sqrt{\text{size}(\mathcal{T})} \|\varphi\|_{1,\infty,\bar{\Omega}}, \quad (38)$$

where  $\sigma = (\sigma_b, \sigma_e)$  (see Definition 3.2). Using inequalities (33), (35)-(38) and formula (10), the equation (31) can be rewritten as

$$- \int_{\Omega} u_{\mathcal{T}}(\mathbf{x}) \Delta \varphi(\mathbf{x}) d\mathbf{x} + \int_{\partial\Omega} \varphi_n(\mathbf{x}) v_{\mathcal{T}}(\mathbf{x}) d\gamma(\mathbf{x}) = \mathbb{T}_3(\varphi, \mathcal{T}), \quad \forall \varphi \in \mathcal{C}^2(\bar{\Omega}), \quad (39)$$

where

$$\mathbb{T}_3(\varphi, \mathcal{T}) = \int_{\Omega} f \varphi_{\mathcal{T}}^{(1)} d\mathbf{x} + \int_{\partial\Omega} g \varphi_{\mathcal{T}}^{(2)} d\gamma(\mathbf{x}) + \alpha \int_{\partial\Omega} \varphi_t v_{\mathcal{T}} d\gamma(\mathbf{x}) + \mathbb{T}_4(\varphi, \mathcal{T}), \quad (40)$$

and

$$|\mathbb{T}_4(\varphi, \mathcal{T})| \leq C_{12} \sqrt{\text{size}(\mathcal{T})}, \quad (41)$$

where  $C_{12}$  depends on  $(\Omega, \alpha, \zeta_1, f, g, \varphi)$ . Writing (39) with  $\mathcal{T} = \mathcal{T}_n$ , using (41) and passing to the limit as  $n$  tends to infinity yield (28).

To prove now that  $u$  satisfies (14), we need the following Lemma which is proven in [EGH 00]:

LEMMA 4.6 ([EGH 00]) Let  $d \geq 1$  and  $\Omega$  be a bounded polygonal open set of  $\mathbb{R}^d$ . Let  $u \in H^1(\Omega)$ ,  $f \in L^2(\Omega)$  and  $v \in L^2(\partial\Omega)$ . We assume that

$$-\int_{\Omega} u(\mathbf{x}) \Delta \varphi(\mathbf{x}) d\mathbf{x} + \int_{\partial\Omega} (\nabla \varphi \cdot \mathbf{n})(\mathbf{x}) v(\mathbf{x}) d\gamma(\mathbf{x}) = \int_{\Omega} f(\mathbf{x}) \varphi(\mathbf{x}) d\mathbf{x}, \quad (42)$$

for every  $\varphi \in \mathcal{C}^2(\overline{\Omega})$  such that  $\varphi = 0$  on  $\partial\Omega$ . Then  $\tilde{\gamma}(u) = v$  a.e. on  $\partial\Omega$ , where  $\tilde{\gamma}(u)$  is the classical trace operator from  $H^1(\Omega)$  to  $L^2(\partial\Omega)$ . (One also has necessarily  $-\Delta u = f$  in  $\mathcal{D}'(\Omega)$ .)

Choosing  $\varphi = 0$  on  $\partial\Omega$  in the equation (28) and using Lemma 4.6, we get  $\tilde{\gamma}(u) = v$  a.e. on  $\partial\Omega$ . Then, an integration by parts in (28) implies that, for any  $\varphi \in \mathcal{D}(\overline{\Omega})$ , we have

$$\begin{aligned} \int_{\Omega} \nabla u(\mathbf{x}) \cdot \nabla \varphi(\mathbf{x}) d\mathbf{x} &= \int_{\Omega} f(\mathbf{x}) \varphi(\mathbf{x}) d\mathbf{x} + \int_{\partial\Omega} g(\mathbf{x}) \varphi(\mathbf{x}) d\gamma(\mathbf{x}) \\ &+ \alpha \int_{\partial\Omega} \varphi_t(\mathbf{x}) \tilde{\gamma}(u)(\mathbf{x}) d\gamma(\mathbf{x}). \end{aligned} \quad (43)$$

Using again an integration by parts and  $\varphi_t = -\varphi_x \mathbf{n}_y + \varphi_y \mathbf{n}_x$ , we obtain

$$\begin{aligned} \int_{\Omega} \nabla u(\mathbf{x}) \cdot \nabla \varphi(\mathbf{x}) d\mathbf{x} &+ \alpha \int_{\Omega} (\varphi_x(\mathbf{x}) u_y(\mathbf{x}) - \varphi_y(\mathbf{x}) u_x(\mathbf{x})) d\mathbf{x} = \int_{\Omega} f(\mathbf{x}) \varphi(\mathbf{x}) d\mathbf{x} \\ &+ \int_{\partial\Omega} g(\mathbf{x}) \varphi(\mathbf{x}) d\gamma(\mathbf{x}), \forall \varphi \in \mathcal{D}(\overline{\Omega}). \end{aligned} \quad (44)$$

Thanks to the density  $\mathcal{D}(\overline{\Omega})$  in  $H^1(\Omega)$ , the formulation (44) is equivalent to (14).

We have proven that the sequence  $u_{\mathcal{T}_n}$  converges to a solution  $u \in H^1(\Omega)$  of (14) in  $L^2(\Omega)$ , up to a subsequence. On the other hand, when the mesh size vanishes in (19), we get (13). Finally, Since the solution  $u$  of (13)-(14) is unique, the whole family  $u_{\mathcal{T}}$  converges to the solution  $u \in H^1(\Omega)$  of (13)-(14) in  $L^2(\Omega)$  and the whole family  $v_{\mathcal{T}}$  converges to  $\tilde{\gamma}(u)$  for the weak topology of  $L^2(\partial\Omega)$  as  $\text{size}(\mathcal{T})$  goes to 0. Now, we prove:

$$\|(u_{\mathcal{T}}, v_{\mathcal{T}})\|_{\star}^2 \rightarrow \int_{\Omega} |\nabla u|^2(\mathbf{x}) d\mathbf{x}, \quad (45)$$

where  $\|(\cdot, \cdot)\|_{\star}$  is the semi-norm defined by

$$\|(u_{\mathcal{T}}, v_{\mathcal{T}})\|_{\star}^2 = |(u_{\mathcal{T}}, v_{\mathcal{T}})|_{1, \mathcal{X}(\mathcal{T})}^2 + \frac{\alpha}{2} |v_{\mathcal{T}}|_{\mathcal{Z}(\mathcal{T})}^2, \quad (46)$$

and the semi-norm  $|(\cdot, \cdot)|_{1, \mathcal{X}(\mathcal{T})}$  (resp.  $|\cdot|_{\mathcal{Z}(\mathcal{T})}$ ) is defined in (7) (resp. (22)). Taking  $v = u$  in (14) leads to

$$\int_{\Omega} |\nabla u|^2(\mathbf{x}) d\mathbf{x} = \int_{\Omega} f(\mathbf{x}) u(\mathbf{x}) d\mathbf{x} + \int_{\partial\Omega} g(\mathbf{x}) \tilde{\gamma}(u) d\gamma(\mathbf{x}). \quad (47)$$

As  $\text{size}(\mathcal{T})$  goes to 0 in (23), we get

$$\|(u_{\mathcal{T}}, v_{\mathcal{T}})\|_{\star}^2 \rightarrow \int_{\Omega} f(\mathbf{x}) u(\mathbf{x}) d\mathbf{x} + \int_{\partial\Omega} g(\mathbf{x}) \tilde{\gamma}(u)(\mathbf{x}) d\gamma(\mathbf{x}). \quad (48)$$

Combining (47) and (48) yields (45).

We have proven the following Theorem:

**THEOREM 4.7 (CONVERGENCE RESULT WHEN  $\alpha$  IS CONSTANT)** Let  $\alpha \in \mathbb{R}$ ,  $\alpha > 0$  (for the other cases, see Remark 3), and  $\zeta_1 > 0$ . Assume Assumption 4.1. Let  $\mathcal{T}$  be an admissible mesh in the sense of Definition 3.1, satisfying the condition (34), and  $\{(f_K, g_\sigma), (K, \sigma) \in \mathcal{T} \times \mathcal{E}_{\text{ext}}\}$  be defined by (9). Let  $(u_{\mathcal{T}}, v_{\mathcal{T}}) \in \mathcal{X}(\mathcal{T})$  be the unique solution of (15)-(19) in the sense of Definition 4.3 (see Theorem 4.4) and let  $u \in H^1(\Omega)$  be the unique solution of (13)-(14). Then:

$$u_{\mathcal{T}} \rightarrow u \text{ in } L^2(\Omega) \text{ as } \text{size}(\mathcal{T}) \rightarrow 0, \quad (49)$$

$$\|(u_{\mathcal{T}}, v_{\mathcal{T}})\|_{\star}^2 \rightarrow \int_{\Omega} |\nabla u|^2(\mathbf{x}) d\mathbf{x}, \text{ as } \text{size}(\mathcal{T}) \rightarrow 0, \quad (50)$$

$$v_{\mathcal{T}} \rightarrow \tilde{\gamma}(u) \text{ in } L^2(\partial\Omega) \text{ for the weak topology as the } \text{size}(\mathcal{T}) \rightarrow 0, \quad (51)$$

where  $\|(\cdot, \cdot)\|_{\star}^2$  is defined by (46) and  $\tilde{\gamma}$  is the classical trace operator from  $H^1(\Omega)$  to  $L^2(\partial\Omega)$ .

#### 4.4 Error estimate

In Section 4.3, we proved the convergence of  $(u_{\mathcal{T}}, v_{\mathcal{T}})$  to the solution  $u$  of (13)-(14) (which belongs to  $H^1(\Omega)$ ). Our aim in this section is to give a convergence order. To do so, we assume that a solution  $u$  of (13)-(14) belongs to  $\mathcal{C}^2(\overline{\Omega})$ . The idea we want to present is mainly based on the idea of the Proof of Theorem 10.1 of [EGH 00] for Neumann Problem. Under the hypotheses of Section 4.3, let  $(u_{\mathcal{T}}, v_{\mathcal{T}}) \in \mathcal{X}(\mathcal{T})$  be the solution of (15)-(19) and consider  $C_{\mathcal{T}} \in \mathbb{R}$  such that

$$\sum_{K \in \mathcal{T}} m(K) \bar{u}(\mathbf{x}_K) = 0, \quad (52)$$

where  $\bar{u} = u + C_{\mathcal{T}}$ . For each  $(K, \sigma) \in \mathcal{T} \times \mathcal{E}_{\text{ext}}$ , let  $e_K = \bar{u}(\mathbf{x}_K) - u_K$  and  $e_{\sigma} = \bar{u}(\mathbf{y}_{\sigma}) - u_{\sigma}$ , where  $u_K$  (resp.  $u_{\sigma}$ ) is the value of  $u_{\mathcal{T}}$  (resp.  $v_{\mathcal{T}}$ ) on  $K$  (resp.  $\sigma$ ) (recall that  $\mathbf{y}_{\sigma}$  is defined in Definition 3.1). We consider  $(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}}) \in \mathcal{X}(\mathcal{T})$  defined by  $(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}}) = (e_K, e_{\sigma})$ , on  $(K, \sigma) \in \mathcal{T} \times \mathcal{E}_{\text{ext}}$ . One defines:

$$R_{K,\sigma} = \frac{\bar{u}(\mathbf{x}_L) - \bar{u}(\mathbf{x}_K)}{d_{K|L}} - \frac{1}{m(\sigma)} \int_{\sigma} \nabla \bar{u}(\mathbf{x}) \cdot \mathbf{n}_{K,\sigma}(\mathbf{x}) d\gamma(\mathbf{x}), \quad \forall \sigma \in \mathcal{E}_{\text{int}} \text{ and } \sigma = K|L, \quad (53)$$

and

$$R_{K,\sigma} = \frac{\bar{u}(\mathbf{y}_{\sigma}) - \bar{u}(\mathbf{x}_K)}{d_{\sigma}} - \frac{1}{m(\sigma)} \int_{\sigma} \nabla \bar{u}(\mathbf{x}) \cdot \mathbf{n}_{K,\sigma}(\mathbf{x}) d\gamma(\mathbf{x}), \quad \forall \sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}. \quad (54)$$

Thus, for  $u \in \mathcal{C}^2(\overline{\Omega})$

$$|R_{K,\sigma}| \leq C_{13} \text{size}(\mathcal{T}), \quad \forall \sigma \in \mathcal{E}_K \text{ and for any } K \in \mathcal{T}, \quad (55)$$

where  $C_{13}$  only depends on  $u$ . Since  $-\Delta \bar{u} = f$ , integrating this equation over any control volume  $K \in \mathcal{T}$  yields:

$$-\int_{\partial K} \nabla \bar{u}(\mathbf{x}) \cdot \mathbf{n}_K d\gamma(\mathbf{x}) = \int_K f(\mathbf{x}) d\mathbf{x}. \quad (56)$$

(Recall that  $\mathbf{n}_K$  is the normal to the boundary  $\partial K$ , outward to  $K$ .) Combining now (53), (54) and (56) we obtain

$$\begin{aligned} & - \sum_{\substack{\sigma \in \mathcal{E}_K \\ \sigma = K|L}} \tau_{K|L} (\bar{u}(\mathbf{x}_L) - \bar{u}(\mathbf{x}_K)) - \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} \tau_{\sigma} (\bar{u}(\mathbf{y}_{\sigma}) - \bar{u}(\mathbf{x}_K)) \\ & = m(K) f_K - \sum_{\sigma \in \mathcal{E}_K} m(\sigma) R_{K,\sigma}, \quad \forall K \in \mathcal{T}. \end{aligned} \quad (57)$$

On the other hand (since  $\nabla \bar{u}(\mathbf{x}) \cdot \mathbf{n}(\mathbf{x}) + \alpha \nabla \bar{u}(\mathbf{x}) \cdot \mathbf{t}(\mathbf{x}) = g(\mathbf{x})$ ,  $\mathbf{x} \in \partial\Omega$ ):

$$\int_{\sigma} \nabla \bar{u}(\mathbf{x}) \cdot \mathbf{n}(\mathbf{x}) d\gamma(\mathbf{x}) + \alpha \int_{\sigma} \nabla \bar{u}(\mathbf{x}) \cdot \mathbf{t}(\mathbf{x}) d\gamma(\mathbf{x}) = m(\sigma)g_{\sigma}, \quad \forall \sigma \in \mathcal{E}_{\text{ext}}, \quad (58)$$

(Recall that  $f_K$  and  $g_{\sigma}$  are defined in (9).) Using (54), formula (10) and (58), we get

$$\tau_{\sigma}(\bar{u}(\mathbf{y}_{\sigma}) - \bar{u}(\mathbf{x}_K)) + \alpha(\bar{u}(\sigma_e) - \bar{u}(\sigma_b)) = m(\sigma)g_{\sigma} + m(\sigma)R_{K,\sigma}, \quad \forall \sigma \in \mathcal{E}_{\text{ext}} \cap \mathcal{E}_K, \quad (59)$$

where  $\sigma = (\sigma_b, \sigma_e)$ , with  $m(\sigma)\mathbf{t} = \sigma_e - \sigma_b$  (see Definition 3.2). This yields

$$\tau_{\sigma}(\bar{u}(\mathbf{y}_{\sigma}) - \bar{u}(\mathbf{x}_K)) + \alpha(\bar{u}(\mathbf{y}_{\sigma}) - \bar{u}(\mathbf{y}_{\sigma-})) + r_{\sigma} - r_{\sigma-} = m(\sigma)g_{\sigma} + m(\sigma)R_{K,\sigma}, \quad (60)$$

where

$$r_{\sigma} = \alpha\{\bar{u}(\sigma_e) - \bar{u}(\mathbf{y}_{\sigma})\}. \quad (61)$$

We have the estimate

$$|r_{\sigma}| \leq C_{14} \text{size}(\mathcal{T}), \quad (62)$$

where  $C_{14}$  only depends on  $(u, \alpha)$ .

Subtracting (15) from (57) and (18) from (60), we get

$$- \sum_{\substack{\sigma \in \mathcal{E}_K \\ \sigma = K|L}} \tau_{\sigma}(e_L - e_K) - \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} \tau_{\sigma}(e_{\sigma} - e_K) = - \sum_{\sigma \in \mathcal{E}_K} m(\sigma)R_{K,\sigma}, \quad \forall K \in \mathcal{T} \quad (63)$$

and

$$\tau_{\sigma}(e_{\sigma} - e_K) = -\alpha(e_{\sigma} - e_{\sigma-}) - r_{\sigma} + r_{\sigma-} + m(\sigma)R_{K,\sigma}, \quad \forall \sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}. \quad (64)$$

Furthermore, subtracting (19) from (52) to get

$$\int_{\Omega} e_{\mathcal{T}}(\mathbf{x}) d\mathbf{x} = 0. \quad (65)$$

Multiplying both sides of equation (63) by  $e_K$ ,  $K \in \mathcal{T}$ , summing over  $K$ ,  $K \in \mathcal{T}$  and using equalities (64) and (20), we get

$$|(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}})|_{1,\mathcal{X}(\mathcal{T})}^2 + \frac{\alpha}{2} |\bar{e}_{\mathcal{T}}|_{\mathcal{Z}(\mathcal{T})}^2 = \mathbb{T}_5^{\mathcal{T}} + \mathbb{T}_6^{\mathcal{T}}, \quad (66)$$

where

$$\mathbb{T}_5^{\mathcal{T}} = - \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} m(\sigma)R_{K,\sigma}e_K, \quad (67)$$

and

$$\mathbb{T}_6^{\mathcal{T}} = \sum_{\sigma \in \mathcal{E}_{\text{ext}}} (-r_{\sigma} + r_{\sigma-})e_{\sigma} + \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} m(\sigma)R_{K,\sigma}e_{\sigma}. \quad (68)$$

We begin with the term  $\mathbb{T}_5^{\mathcal{T}}$ , reordering the sum, using the fact that  $R_{K,\sigma} = -R_{L,\sigma}$ , for all  $\sigma \in \mathcal{E}_{\text{int}}$  and  $\sigma = K|L$ , and Inequality (55), we get

$$\begin{aligned} |\mathbb{T}_5^{\mathcal{T}}| &= \left| \sum_{\substack{\sigma \in \mathcal{E}_{\text{int}} \\ \sigma = K|L}} m(\sigma)R_{K,\sigma}(e_L - e_K) + \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} m(\sigma)R_{K,\sigma}e_K \right| \\ &\leq C_{15} \text{size}(\mathcal{T}) \left( |e_{\mathcal{T}}|_{1,\mathcal{T}} + \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} m(\sigma)|e_K| \right) \\ &\leq C_{15} \text{size}(\mathcal{T}) \left( |(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}})|_{1,\mathcal{X}(\mathcal{T})} + \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} m(\sigma)|e_K| \right), \end{aligned} \quad (69)$$

where  $C_{15}$  only depends on  $(u, \Omega)$ . In order to estimate the second term on the r.h.s. (right hand side) of (69), we consider the discrete trace  $\bar{\gamma}(e_{\mathcal{T}})$  of  $e_{\mathcal{T}}$  (see [EGH 00], page 807) which is defined by  $\bar{\gamma}(e_{\mathcal{T}}) = e_K$  a.e. (for the one-dimensional Lebesgue measure) on  $\sigma$ , if  $\sigma \in \mathcal{E}_{\text{ext}} \cap \mathcal{E}_K$ . Using Cauchy-Schwarz inequality and inequalities (10.10) and (10.25) of [EGH 00] with the fact that  $e_{\mathcal{T}}$  satisfies (65), we get:

$$\begin{aligned} \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} m(\sigma) |e_K| &\leq C_{16} \|\bar{\gamma}(e_{\mathcal{T}})\|_{L^2(\partial\Omega)} \\ &\leq C_{17} (|e_{\mathcal{T}}|_{1,\mathcal{T}} + \|e_{\mathcal{T}}\|_{L^2(\Omega)}) \\ &\leq C_{18} |e_{\mathcal{T}}|_{1,\mathcal{T}}, \end{aligned} \quad (70)$$

where  $C_{16}, C_{17}$  and  $C_{18}$  are only depending on  $\Omega$ . With (69), this yields:

$$|\mathbb{T}_5^{\mathcal{T}}| \leq C_{19} \text{size}(\mathcal{T}) \left( |(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}})|_{1,\mathcal{X}(\mathcal{T})}^2 + \frac{\alpha}{2} |\bar{e}_{\mathcal{T}}|_{\mathcal{Z}(\mathcal{T})}^2 \right)^{\frac{1}{2}}, \quad (71)$$

where  $C_{19} = C_{15}(1 + C_{18})$ .

To estimate  $\mathbb{T}_6^{\mathcal{T}}$ , we have to assume that the mesh  $\mathcal{T}$  satisfies, for some  $\zeta_2 > 0$ :

$$m(\sigma) \geq \zeta_2 \text{size}(\mathcal{T}), \quad \forall \sigma \in \mathcal{E}_{\text{ext}}. \quad (72)$$

*Remark 5* The condition (72) implies the condition (34) with  $\zeta_1 = \frac{2}{\zeta_2}$ .

Using inequalities (55) and (62) and using the the assumption (72), we get

$$\begin{aligned} |\mathbb{T}_6^{\mathcal{T}}| &\leq \left| \sum_{\sigma \in \mathcal{E}_{\text{ext}}} (-r_{\sigma} + r_{\sigma-}) e_{\sigma} \right| + \left| \sum_{\sigma \in \mathcal{E}_{\text{ext}}} m(\sigma) R_{K,\sigma} e_{\sigma} \right| \\ &\leq \left| \sum_{\sigma \in \mathcal{E}_{\text{ext}}} r_{\sigma} (e_{\sigma+} - e_{\sigma}) \right| + C_{20} \text{size}(\mathcal{T}) \left( \sum_{\sigma \in \mathcal{E}_{\text{ext}}} m(\sigma) e_{\sigma}^2 \right)^{\frac{1}{2}} \\ &\leq C_{14} \text{size}(\mathcal{T}) \left( \sum_{\sigma \in \mathcal{E}_{\text{ext}}} 1 \right)^{\frac{1}{2}} |\bar{e}_{\mathcal{T}}|_{\mathcal{Z}(\mathcal{T})} + C_{20} \text{size}(\mathcal{T}) \left( \sum_{\sigma \in \mathcal{E}_{\text{ext}}} m(\sigma) e_{\sigma}^2 \right)^{\frac{1}{2}} \\ &\leq C_{21} \sqrt{\text{size}(\mathcal{T})} |\bar{e}_{\mathcal{T}}|_{\mathcal{Z}(\mathcal{T})} + C_{20} \text{size}(\mathcal{T}) \left( \sum_{\sigma \in \mathcal{E}_{\text{ext}}} m(\sigma) e_{\sigma}^2 \right)^{\frac{1}{2}}, \end{aligned} \quad (73)$$

where  $C_{20}$  depends on  $(u, \Omega)$  and  $C_{21}$  depends on  $(u, \Omega, \zeta_2, \alpha)$ . To estimate the second term on the r.h.s. of (73), we use the triangular inequality and inequalities (10.10) and (10.25) of [EGH 00] with the fact that  $e_{\mathcal{T}}$  satisfies (65) (and that  $d_{\sigma} < \text{diam}(\Omega)$ ). We get:

$$\begin{aligned} \left( \sum_{\sigma \in \mathcal{E}_{\text{ext}}} m(\sigma) e_{\sigma}^2 \right)^{\frac{1}{2}} &\leq \left( \sum_{\sigma \in \mathcal{E}_{\text{ext}}} m(\sigma) (e_{\sigma} - e_K)^2 \right)^{\frac{1}{2}} + \left( \sum_{\sigma \in \mathcal{E}_{\text{ext}}} m(\sigma) e_K^2 \right)^{\frac{1}{2}} \\ &\leq C_{22} |(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}})|_{1,\mathcal{X}(\mathcal{T})}, \end{aligned} \quad (74)$$

where  $C_{22}$  only depends on  $\Omega$ . Combining inequalities (73) and (74) gives:

$$|\mathbb{T}_6^{\mathcal{T}}| \leq C_{23} \sqrt{\text{size}(\mathcal{T})} \left( |(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}})|_{1,\mathcal{X}(\mathcal{T})}^2 + \frac{\alpha}{2} |\bar{e}_{\mathcal{T}}|_{\mathcal{Z}(\mathcal{T})}^2 \right)^{\frac{1}{2}}, \quad (75)$$

where  $C_{23}$  only depends on  $(u, \Omega, \zeta_2, \alpha)$ . Equality (66) with inequalities (71), (75) implies:

$$\left( |(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}})|_{1,\mathcal{X}(\mathcal{T})}^2 + \frac{\alpha}{2} |\bar{e}_{\mathcal{T}}|_{\mathcal{Z}(\mathcal{T})}^2 \right)^{\frac{1}{2}} \leq C_{24} \sqrt{\text{size}(\mathcal{T})}, \quad (76)$$

where  $C_{24}$  only depends on  $(u, \Omega, \zeta_2, \alpha)$ . Inequality (76) together with (65) and the discrete mean Poincaré inequality (10.10) of [EGH 00] yields:

$$\|e_{\mathcal{T}}\|_{L^2(\Omega)}^2 \leq C_{25} \text{size}(\mathcal{T}), \quad (77)$$

where  $C_{25}$  only depends on  $(u, \Omega, \zeta_2, \alpha)$ .

In order to obtain an error estimate between  $u$  and  $u_{\mathcal{T}}$  (instead of  $\bar{u}$  and  $u_{\mathcal{T}}$ ), we first give an estimate on  $C_{\mathcal{T}}$ . Using the fact that  $\int_{\Omega} u(\mathbf{x}) \, d\mathbf{x} = 0$  (Equation (13)) and Equation (52), we get

$$\begin{aligned} m(\Omega)C_{\mathcal{T}} &= \int_{\Omega} \{u(\mathbf{x}) - \bar{u}(\mathbf{x})\} d\mathbf{x} = \sum_{K \in \mathcal{T}} \int_K \bar{u}(\mathbf{x}) d\mathbf{x} \\ &= \sum_{K \in \mathcal{T}} \int_K \{\bar{u}(\mathbf{x}_K) + \nabla u(\psi(\mathbf{x})) \cdot (\mathbf{x} - \mathbf{x}_K)\} d\mathbf{x} = \sum_{K \in \mathcal{T}} \int_K \nabla u(\psi(\mathbf{x})) \cdot (\mathbf{x} - \mathbf{x}_K) d\mathbf{x}, \end{aligned}$$

and then:

$$m(\Omega)|C_{\mathcal{T}}| \leq C_{26} \text{size}(\mathcal{T}) \sum_{K \in \mathcal{T}} m(K) \leq C_{27} \text{size}(\mathcal{T}), \quad (78)$$

where  $C_{26} = \|\nabla u\|_{L^\infty(\Omega)}$ ,  $C_{27} = C_{26} m(\Omega)$  and  $\psi(\mathbf{x})$  is a some point between the points  $\mathbf{x}$  and  $\mathbf{x}_K$ . Furthermore, we have

$$\begin{aligned} \sum_{K \in \mathcal{T}} \int_K \{u(\mathbf{x}_K) - u(\mathbf{x})\}^2 d\mathbf{x} &\leq (C_{26})^2 (\text{size}(\mathcal{T}))^2 \sum_{K \in \mathcal{T}} m(K) \\ &= (C_{26})^2 m(\Omega) (\text{size}(\mathcal{T}))^2. \end{aligned} \quad (79)$$

Using triangular Inequality and Error Estimate (77) combined with inequalities (78) and (79), we get the following error estimate in  $L^2(\Omega)$ -norm:

$$\begin{aligned} \|u_{\mathcal{T}} - u\|_{L^2(\Omega)}^2 &= \sum_{K \in \mathcal{T}} \int_K \{u_K - u(\mathbf{x})\}^2 d\mathbf{x} \\ &\leq 3\|e_{\mathcal{T}}\|_{L^2(\Omega)}^2 + 3m(\Omega)C_{\mathcal{T}}^2 + 3 \sum_{K \in \mathcal{T}} \int_K \{u(\mathbf{x}_K) - u(\mathbf{x})\}^2 d\mathbf{x} \\ &\leq C_{28} \text{size}(\mathcal{T}), \end{aligned} \quad (80)$$

where  $C_{28}$  depends only on  $(u, \Omega, \zeta_2, \alpha)$ .

We now turn to get an error estimate in a discrete  $H_0^1(\Omega)$  norm. For each  $(K, \sigma) \in \mathcal{T} \times \mathcal{E}_{\text{ext}}$ , let  $(e_K^{\text{real}}, e_\sigma^{\text{real}}) = (u(\mathbf{x}_K) - u_K, u(\mathbf{y}_\sigma) - u_\sigma)$  and consider  $(e_{\mathcal{T}}^{\text{real}}, \bar{e}_{\mathcal{T}}^{\text{real}}) \in \mathcal{X}(\mathcal{T})$  defined by  $(e_{\mathcal{T}}^{\text{real}}, \bar{e}_{\mathcal{T}}^{\text{real}}) = (e_K^{\text{real}}, e_\sigma^{\text{real}})$ , on  $(K, \sigma) \in \mathcal{T} \times \mathcal{E}_{\text{ext}}$ . We remark that (because  $\bar{u} = u + C_{\mathcal{T}}$ )

$$|(e_{\mathcal{T}}^{\text{real}}, \bar{e}_{\mathcal{T}}^{\text{real}})|_{1, \mathcal{X}(\mathcal{T})}^2 + \frac{\alpha}{2} |\bar{e}_{\mathcal{T}}^{\text{real}}|_{\mathcal{Z}(\mathcal{T})}^2 = |(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}})|_{1, \mathcal{X}(\mathcal{T})}^2 + \frac{\alpha}{2} |\bar{e}_{\mathcal{T}}|_{\mathcal{Z}(\mathcal{T})}^2. \quad (81)$$

With (76), this implies

$$\left( |(e_{\mathcal{T}}^{\text{real}}, \bar{e}_{\mathcal{T}}^{\text{real}})|_{1, \mathcal{X}(\mathcal{T})}^2 + \frac{\alpha}{2} |\bar{e}_{\mathcal{T}}^{\text{real}}|_{\mathcal{Z}(\mathcal{T})}^2 \right)^{\frac{1}{2}} \leq C_{24} \sqrt{\text{size}(\mathcal{T})}. \quad (82)$$

We have proven the following Error Estimate:

**THEOREM 4.8 ( $\mathcal{C}^2$ -ERROR ESTIMATE WHEN  $\alpha$  IS CONSTANT)** Let  $\alpha \in \mathbb{R}$ ,  $\alpha > 0$  (for the other cases, see Remark 3), and  $\zeta_2 > 0$ . Assume Assumption 4.1. Let  $\mathcal{T}$  be an admissible



mesh in the sense of Definition 3.1, satisfying the condition (72), and  $\{(f_K, g_\sigma), (K, \sigma) \in \mathcal{T} \times \mathcal{E}_{\text{ext}}\}$  be defined by (9). Let  $(u_{\mathcal{T}}, v_{\mathcal{T}}) \in \mathcal{X}(\mathcal{T})$  be the unique solution of (15)-(19) in the sense of Definition 4.3 (see Theorem 4.4) and let  $u$  be the unique solution of (13)-(14). Assume that  $u \in \mathcal{C}^2(\overline{\Omega})$ . Then, there exist  $(C_{28}, C_{24})$  only depending on  $(u, \Omega, \zeta_2, \alpha)$  such that:

$$\|u_{\mathcal{T}} - u\|_{L^2(\Omega)}^2 \leq C_{28} \text{size}(\mathcal{T}), \quad (83)$$

and

$$\left( |(e_{\mathcal{T}}^{\text{real}}, \bar{e}_{\mathcal{T}}^{\text{real}})|_{1, \mathcal{X}(\mathcal{T})}^2 + \frac{\alpha}{2} |\bar{e}_{\mathcal{T}}^{\text{real}}|_{\mathcal{Z}(\mathcal{T})}^2 \right)^{\frac{1}{2}} \leq C_{24} \sqrt{\text{size}(\mathcal{T})}, \quad (84)$$

where  $(e_{\mathcal{T}}^{\text{real}}, \bar{e}_{\mathcal{T}}^{\text{real}}) \in \mathcal{X}(\mathcal{T})$  (see Definitions 3.3 and 3.4) is defined by  $(e_{\mathcal{T}}^{\text{real}}, \bar{e}_{\mathcal{T}}^{\text{real}}) = (u(\mathbf{x}_K) - u_K, u(\mathbf{y}_\sigma) - u_\sigma)$ , on  $(K, \sigma) \in \mathcal{T} \times \mathcal{E}_{\text{ext}}$ .

Under the hypotheses of Theorem 4.8, one deduces from (84) the following estimate for the  $L^2$  norm of the flux:

$$\begin{aligned} & \sum_{\substack{\sigma \in \mathcal{E}_{\text{int}} \\ \sigma = K|L}} m(\sigma) d_\sigma \left( \frac{u_L - u_K}{d_\sigma} - \frac{1}{m(\sigma)} \int_\sigma \nabla u(\mathbf{x}) \cdot \mathbf{n}_{K, \sigma} d\gamma(\mathbf{x}) \right)^2 \\ & + \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} m(\sigma) d_\sigma \left( \frac{u_\sigma - u_K}{d_\sigma} - \frac{1}{m(\sigma)} \int_\sigma \nabla u(\mathbf{x}) \cdot \mathbf{n}_{K, \sigma} d\gamma(\mathbf{x}) \right)^2 \leq (C_{24})^2 \text{size}(\mathcal{T}). \end{aligned} \quad (85)$$

## 5 The case “ $\alpha$ smooth”

This section is a generalization of Section 4. We now consider Problem (1) when  $\alpha$  is a smooth function, namely  $\alpha \in \mathcal{C}^1(\overline{\Omega})$ . Problem (1) reads:

$$\begin{cases} -\Delta u(\mathbf{x}) = f(\mathbf{x}), & \mathbf{x} \in \Omega, \\ u_n(\mathbf{x}) + (\alpha u)_t(\mathbf{x}) = g(\mathbf{x}), & \mathbf{x} \in \partial\Omega. \end{cases} \quad (86)$$

In order to get (for some functions  $\alpha$ , see Theorem 5.2), the existence of a solution for Problem (86), we assume that:

ASSUMPTION 5.1

- (i)  $(f, g) \in L^2(\Omega) \times L^2(\partial\Omega)$  and  $\int_\Omega f(\mathbf{x}) d\mathbf{x} + \int_{\partial\Omega} g(\mathbf{x}) d\gamma(\mathbf{x}) = 0$ .
- (ii)  $\alpha \in \mathcal{C}^1(\overline{\Omega})$ .

THEOREM 5.2 Under the Assumption 5.1, let  $C_\alpha = \min_{\partial\Omega} \alpha_t$ . Then there exists  $\delta < 0$ , only depending on  $\Omega$ , such that if  $\alpha$  satisfies the condition  $C_\alpha \geq \delta$ , then there exists a unique solution to (87)-(88):

$$u \in H^1(\Omega), \int_\Omega u(\mathbf{x}) d\mathbf{x} = 0, \quad (87)$$

$$b(u, v) = F(v), \forall v \in H^1(\Omega), \quad (88)$$

where

$$b(u, v) = \int_\Omega \nabla u(\mathbf{x}) \cdot \nabla v(\mathbf{x}) d\mathbf{x} + \int_\Omega \{(\alpha u)_y(\mathbf{x}) v_x(\mathbf{x}) - (\alpha u)_x(\mathbf{x}) v_y(\mathbf{x})\} d\mathbf{x}, \quad (89)$$

and

$$F(v) = \int_\Omega f(\mathbf{x}) v(\mathbf{x}) d\mathbf{x} + \int_{\partial\Omega} g(\mathbf{x}) v(\mathbf{x}) d\gamma(\mathbf{x}). \quad (90)$$

*Remark 6* Problem (87)-(88) appears to be a weak formulation of Problem (86) (with the additional assumption  $\int_{\Omega} u(\mathbf{x}) d\mathbf{x} = 0$ , in order to obtain also a uniqueness result) since for  $\alpha u$  and  $v$  regular enough (say, for instance,  $\alpha u \in H^2(\Omega)$  and  $v \in H^1(\Omega)$ ) one has:

$$\int_{\partial\Omega} (\alpha u)_t(\mathbf{x}) v(\mathbf{x}) d\gamma(\mathbf{x}) = \int_{\Omega} \{(\alpha u)_y(\mathbf{x}) v_x(\mathbf{x}) - (\alpha u)_x(\mathbf{x}) v_y(\mathbf{x})\} d\mathbf{x}.$$

**Proof** To prove the existence and uniqueness of the solution of (87)-(88), we apply the classical Lax-Milgram lemma. Thanks to Assumption 5.1, It is clear that  $b(\cdot, \cdot)$  and  $F(\cdot)$  are continuous on  $H^1(\Omega) \times H^1(\Omega)$  and  $H^1(\Omega)$  respectively.

In order to prove the coercivity of  $b$  (under an additional assumption on  $\alpha_t$ ), let  $u \in H^1(\Omega)$ . Using  $u^2 \in W^{1,1}(\Omega)$  and  $\alpha \in C^1(\Omega)$ , one obtains:

$$\begin{aligned} b(u, u) &= \int_{\Omega} |\nabla u|^2(\mathbf{x}) d\mathbf{x} + \int_{\Omega} \{(\alpha u)_y(\mathbf{x}) u_x(\mathbf{x}) - (\alpha u)_x(\mathbf{x}) u_y(\mathbf{x})\} d\mathbf{x} \\ &= \int_{\Omega} |\nabla u|^2(\mathbf{x}) d\mathbf{x} + \frac{1}{2} \int_{\Omega} \{\alpha_y(\mathbf{x}) (u^2)_x(\mathbf{x}) - \alpha_x(\mathbf{x}) (u^2)_y(\mathbf{x})\} d\mathbf{x} \\ &= \int_{\Omega} |\nabla u|^2(\mathbf{x}) d\mathbf{x} + \frac{1}{2} \int_{\partial\Omega} \alpha_t(\mathbf{x}) u^2(\mathbf{x}) d\gamma(\mathbf{x}). \end{aligned} \quad (91)$$

Then, with  $C_{\alpha} = \min_{\partial\Omega} \alpha_t$  (note that  $C_{\alpha} \leq 0$  since the mean value of  $\alpha_t$  on  $\partial\Omega$  is 0):

$$b(u, u) \geq |u|_{1,\Omega}^2 + \frac{C_{\alpha}}{2} \int_{\partial\Omega} u^2(\mathbf{x}) d\gamma(\mathbf{x}).$$

This gives, using the continuity of the trace operator from  $H^1(\Omega)$  to  $L^2(\partial\Omega)$ ,

$$b(u, u) \geq |u|_{1,\Omega}^2 + \frac{C_{\alpha} C_{29}}{2} |u|_{1,\Omega}^2 + \frac{C_{\alpha} C_{29}}{2} \|u\|_{0,\Omega}^2, \quad (92)$$

where  $C_{29}$  is a positive number only depending on  $\Omega$ . Let

$$\mathcal{H} = \{v \in H^1(\Omega), \int_{\Omega} v(\mathbf{x}) d\mathbf{x} = 0\}.$$

Inequality (92) with the mean Poincaré inequality implies:

$$b(u, u) \geq |u|_{1,\Omega}^2 + \frac{C_{\alpha} C_{29}}{2} |u|_{1,\Omega}^2 + \frac{C_{\alpha} C_{30} C_{29}}{2} |u|_{1,\Omega}^2, \quad \forall u \in \mathcal{H},$$

where  $C_{30}$  is a positive number only depending on  $\Omega$ . The previous inequality can be written as:

$$b(u, u) \geq (1 + C_{31} C_{\alpha}) |u|_{1,\Omega}^2, \quad \forall u \in \mathcal{H},$$

where

$$C_{31} = \frac{(1 + C_{30}) C_{29}}{2}. \quad (93)$$

Thus for

$$C_{\alpha} > -\frac{1}{C_{31}} = \delta,$$

the bilinear form  $b(\cdot, \cdot)$  becomes coercive on  $\mathcal{H}$ . Then, by the Lax-Milgram lemma, there exists a unique  $u \in \mathcal{H}$  such that

$$b(u, v) = F(v), \quad \forall v \in \mathcal{H}. \quad (94)$$

Consider now  $v \in H^1(\Omega)$  and  $C \in \mathbb{R}$  such that

$$\int_{\Omega} \bar{v}(\mathbf{x}) d\mathbf{x} = 0,$$

where  $\bar{v} = v + C$ .

Since  $b(u, C) = 0$ ; using the item (i) of Assumption 5.1 and (94), we get

$$b(u, v) = F(v), \forall v \in H^1(\Omega), \quad (95)$$

which completes the Proof. ■

## 5.1 The finite volume scheme for (86)

To present the finite volume scheme for (86) and to analyze its convergence, we need some more notations and an additional semi-norm (we also use the semi-norm of Definition 3.4).

**DEFINITION 5.3** Let  $\alpha \in \mathcal{C}^1(\overline{\Omega})$ . For  $\sigma \in \mathcal{E}_{\text{ext}}$ , the notations  $\sigma_e$ ,  $\sigma_b$ ,  $u_{\sigma+}$  and  $u_{\sigma-}$  are given in Definition 3.2 (in particular,  $\sigma = (\sigma_b, \sigma_e)$ , with  $\mathbf{m}(\sigma)\mathbf{t} = \sigma_e - \sigma_b$ ). We set:

$$\begin{aligned} u_{\sigma,+} &= u_{\sigma} \text{ and } u_{\sigma,-} = u_{\sigma+} \text{ if } \alpha(\sigma_e) \geq 0, \\ u_{\sigma,+} &= u_{\sigma+} \text{ and } u_{\sigma,-} = u_{\sigma} \text{ if } \alpha(\sigma_e) < 0. \end{aligned}$$

Let  $\mathcal{Z}(\mathcal{T})$  be the space of functions which are constant on each  $\sigma \in \mathcal{E}_{\text{ext}}$  (see Definition 3.3). We define the following semi-norm on  $\mathcal{Z}(\mathcal{T})$ :

$$|v_{\mathcal{T}}|_{\alpha, \mathcal{Z}(\mathcal{T})}^2 = \sum_{\sigma \in \mathcal{E}_{\text{ext}}} |\alpha(\sigma_e)| (u_{\sigma,+} - u_{\sigma,-})^2.$$

*Remark 7* Since  $\{u_{\sigma,+}, u_{\sigma,-}\} = \{u_{\sigma}, u_{\sigma+}\}$ , for all  $\sigma \in \mathcal{E}_{\text{ext}}$ , the semi-norm  $|\cdot|_{\alpha, \mathcal{Z}(\mathcal{T})}$  can be written as:

$$|v_{\mathcal{T}}|_{\alpha, \mathcal{Z}(\mathcal{T})}^2 = \sum_{\sigma \in \mathcal{E}_{\text{ext}}} |\alpha(\sigma_e)| (u_{\sigma+} - u_{\sigma})^2.$$

*Remark 8* Let  $\sigma = (\sigma_b, \sigma_e) \in \mathcal{E}_{\text{ext}}$ . The value  $u_{\sigma,+}$  is the upstream (w.r.t. the sign of  $\alpha(\sigma_e)$ ) value of  $u$  at point  $\sigma_e$ . In order to get the stability of the scheme, it will be used to discretize the tangential term  $(\alpha u)_t$  (see equation (97)). Such an upstream choice is classical to discretize the convection term of elliptic problems in two or three Dimensions (see, for instance, [EGH 00] pages 766 and 767).

Let  $(u_K)_{K \in \mathcal{T}}$  and  $(u_{\sigma})_{\sigma \in \mathcal{E}_{\text{ext}}}$  denote the discrete unknowns. The numerical scheme for the discretization of (86) is given by the following set of equations:

$$\sum_{\mathcal{E}_K} F_{K,\sigma} = \mathbf{m}(K) f_K, \forall K \in \mathcal{T}, \quad (96)$$

where  $F_{K,\sigma}$  is defined in (16)-(17)),

$$\tau_{\sigma}(u_{\sigma} - u_K) = -\alpha(\sigma_e)u_{\sigma,+} + \alpha(\sigma_b)u_{\sigma-,+} + \mathbf{m}(\sigma)g_{\sigma}, \forall \sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}, \quad (97)$$

where the notation  $u_{\sigma,+}$  is defined in Definition 5.3 and  $\sigma = (\sigma_b, \sigma_e)$  with  $\mathbf{m}(\sigma)\mathbf{t} = \sigma_e - \sigma_b$  (note also that  $\sigma_b = (\sigma^-)_e$ ), and, for the discretization of the equation (87),

$$\sum_{K \in \mathcal{T}} \mathbf{m}(K) u_K = 0. \quad (98)$$

## 5.2 Existence and uniqueness of the finite volume solution

In this Section, we prove the existence and the uniqueness of the solution to (96)-(98), under a condition on  $\alpha$  as in Theorem 5.2. A first method to prove this existence and uniqueness result is to use the Lax-Milgram lemma (as in Theorem 5.2). We present here another method, similar to that given for proving Theorem 4.4, but slightly more complicated since the constant vectors are not solutions of the homogeneous system associated to (96)-(97) (with  $F_{K,\sigma}$  given by (16)-(17)).

Assume Assumption 5.1 and let  $\mathcal{T}$  be an admissible mesh in the sense of Definition 3.1. Let  $M_1$  be the number of elements of  $\mathcal{T}$ ,  $M_2$  the number of elements of  $\mathcal{E}_{\text{ext}}$  and  $M = M_1 + M_2$ . The system (96)-(98) (with  $F_{K,\sigma}$  given by (16)-(17)) can be viewed as a system of  $M$  unknowns (which are  $\{u_K, K \in \mathcal{T}\}$  and  $\{u_\sigma, \sigma \in \mathcal{E}_{\text{ext}}\}$ ) with  $M+1$  equations. After the choice of an order for the unknowns and the equations, it can be written as  $Aw = b$ , where  $A$  is  $(M+1) \times M$  matrix,  $w \in \mathbb{R}^M$  is the unknown vector and  $b \in \mathbb{R}^{M+1}$  is given by the data (namely  $f$  and  $g$ ). We first prove that the null space of this matrix  $A$  is reduced to the null vector.

Indeed, let  $\{u_K\}_K$  and  $\{u_\sigma\}_\sigma$  be a solution of (96)-(98) (with (16)-(17)) with  $(f_K, g_\sigma) = (0, 0)$  for all  $(K, \sigma) \in \mathcal{T} \times \mathcal{E}_{\text{ext}}$ . Multiplying both sides of equation (96) by  $u_K$ ,  $K \in \mathcal{T}$ , summing over  $K \in \mathcal{T}$  and using equation (97) yields

$$|(u_{\mathcal{T}}, v_{\mathcal{T}})|_{1, \mathcal{X}(\mathcal{T})}^2 + \sum_{\sigma \in \mathcal{E}_{\text{ext}}} \{\alpha(\sigma_e)u_{\sigma,+} - \alpha(\sigma_b)u_{\sigma,-,+}\}u_\sigma = 0, \quad (99)$$

where  $(u_{\mathcal{T}}, v_{\mathcal{T}}) \in \mathcal{X}(\mathcal{T})$  (see Definition 3.3) and  $u_{\mathcal{T}}(\mathbf{x}) = u_K$ , for  $\mathbf{x} \in K$ , and  $v_{\mathcal{T}}(\mathbf{x}) = u_\sigma$ , for  $\mathbf{x} \in \sigma$ , for all  $(K, \sigma) \in \mathcal{T} \times \mathcal{E}_{\text{ext}}$ . The second term on the l.h.s. of (99) can be written as follows

$$\sum_{\sigma \in \mathcal{E}_{\text{ext}}} \{\alpha(\sigma_e)u_{\sigma,+} - \alpha(\sigma_b)u_{\sigma,-,+}\}u_\sigma = \sum_{\sigma \in \mathcal{E}_{\text{ext}}} |\alpha(\sigma_e)| \{u_{\sigma,+} - u_{\sigma,-}\}u_{\sigma,+}. \quad (100)$$

(recall that  $u_{\sigma,-}$  is defined in Definition 5.3.) This gives that (similar techniques are used in [EGH 00], page 769)

$$\sum_{\sigma \in \mathcal{E}_{\text{ext}}} \{\alpha(\sigma_e)u_{\sigma,+} - \alpha(\sigma_b)u_{\sigma,-,+}\}u_\sigma = \frac{1}{2} \left( \sum_{\sigma \in \mathcal{E}_{\text{ext}}} |\alpha(\sigma_e)| \{ (u_{\sigma,+} - u_{\sigma,-})^2 + u_{\sigma,+}^2 - u_{\sigma,-}^2 \} \right). \quad (101)$$

The second term on the r.h.s. of equality (101) can be written as

$$\sum_{\sigma \in \mathcal{E}_{\text{ext}}} |\alpha(\sigma_e)| \{u_{\sigma,+}^2 - u_{\sigma,-}^2\} = \int_{\partial\Omega} \alpha_t(\mathbf{x}) v_{\mathcal{T}}^2(\mathbf{x}) d\gamma(\mathbf{x}). \quad (102)$$

Combining now (99)-(102), we get

$$|(u_{\mathcal{T}}, v_{\mathcal{T}})|_{1, \mathcal{X}(\mathcal{T})}^2 + \frac{1}{2} |v_{\mathcal{T}}|_{\alpha, \mathcal{Z}(\mathcal{T})}^2 + \frac{1}{2} \int_{\partial\Omega} \alpha_t(\mathbf{x}) v_{\mathcal{T}}^2(\mathbf{x}) d\gamma(\mathbf{x}) = 0. \quad (103)$$

We prove now that, using (98) and under some condition on  $C_\alpha = \min\{\alpha_t(\mathbf{x}), \mathbf{x} \in \partial\Omega\}$  (this condition is similar to that of Theorem 5.2), the l.h.s. of (103) vanishes if and only if  $u_K = u_\sigma = 0$ , for all  $(K, \sigma) \in \mathcal{T} \times \mathcal{E}_{\text{ext}}$ . Indeed, using Inequality  $(a+b)^2 \leq 2a^2 + 2b^2$ , we have:

$$\int_{\partial\Omega} \alpha_t(\mathbf{x}) v_{\mathcal{T}}^2(\mathbf{x}) d\gamma(\mathbf{x}) \geq 2C_\alpha \left( \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} m(\sigma)(u_\sigma - u_K)^2 + \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} m(\sigma)u_K^2 \right). \quad (104)$$

Combining inequalities (10.10) and (10.25) of [EGH 00] with Equation (98) implies that

$$\sum_{\sigma \in \mathcal{E}_{\text{ext}}} m(\sigma) u_K^2 \leq C_{32} |u_{\mathcal{T}}|_{1,\mathcal{T}}^2, \quad (105)$$

where  $C_{32}$  is a constant only depending on  $\Omega$ . Equality (103) with inequalities (104) and (105) implies

$$(1 + C_\alpha C_{32}) |u_{\mathcal{T}}|_{1,\mathcal{T}}^2 + \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} \frac{m(\sigma)}{d_\sigma} (1 + d_\sigma C_\alpha) (u_K - u_\sigma)^2 + \frac{1}{2} |v_{\mathcal{T}}|_{\alpha,\mathcal{Z}(\mathcal{T})}^2 \leq 0. \quad (106)$$

Let  $\delta = \min\{-1/C_{32}, -1/\text{diam}(\Omega)\}$  and assume  $C_\alpha > \delta$  (note that necessarily  $C_\alpha \leq 0$ ). Then, all the terms on the left hand side of (106) are positive (since, in particular,  $d_\sigma \leq \text{diam}(\Omega)$  for all  $\sigma \in \mathcal{E}_{\text{ext}}$ ). This gives that there exists  $C \in \mathbb{R}$  such that  $u_K = C$  for all  $K \in \mathcal{T}$  and  $u_\sigma = C$  for all  $\sigma \in \mathcal{E}_{\text{ext}}$  (as in Remark 1). Finally, thanks to (98), one deduces  $C = 0$ . This proves that the null space of the matrix  $A$  is reduced to the null vector.

*Remark 9* Indeed, assuming only  $C_\alpha > -1/C_{32}$ , if  $\text{size}(\mathcal{T})$  is small enough (more precisely, if  $\text{size}(\mathcal{T}) < -1/C_\alpha$ ), all the terms on the left hand side of (106) are positive. Then, we can also conclude that the null space of the matrix  $A$  is reduced to the null vector. Therefore, as in Theorem 5.4 below, assuming Assumption 5.1, we have existence and uniqueness of the solution of (96)-(98) (with (16)-(17)), when  $\mathcal{T}$  be an admissible mesh in the sense of Definition 3.1.

Since the dimension of the null space of  $A$  is 0, the dimension of the range of  $A$  is  $M$ . But the range of  $A$  is included in  $\mathbb{R}^{M+1}$  and the following condition is necessary for (96)-(98) (with (16)-(17)) to have a solution:

$$\sum_{K \in \mathcal{T}} m(K) f_K + \sum_{\sigma \in \mathcal{E}_{\text{ext}}} m(\sigma) g_\sigma = 0. \quad (107)$$

Then, Condition (107) is also a sufficient condition for (96)-(98) (with (16)-(17)) to have a solution. Finally, under Condition (107) (which is a consequence of Assumption 5.1) the system (96)-(98) (with (16)-(17)) has a unique solution. (By the way, the same result of existence and uniqueness remains true if the right hand side of Equation (98) is replaced by any real value.)

This proves the following Theorem:

**THEOREM 5.4** Assume Assumption 5.1 and let  $C_\alpha = \min_{\partial\Omega} \alpha_t$ . Let  $\mathcal{T}$  be an admissible mesh in the sense of Definition 3.1. Then, There exists  $\delta < 0$ , only depending on  $\Omega$ , such that if  $C_\alpha > \delta$  the system (96)-(98) (with (16)-(17)), where  $\{(f_K, g_\sigma) : (K, \sigma) \in \mathcal{T} \times \mathcal{E}_{\text{ext}}\}$  is given by (9), has a unique solution.

Theorem 5.4 allows us to introduce the following Definition.

**DEFINITION 5.5** An element  $(u_{\mathcal{T}}, v_{\mathcal{T}}) \in \mathcal{X}(\mathcal{T})$  (see Definition 3.3) is the solution of (96)-(98) if  $u_{\mathcal{T}}(\mathbf{x}) = u_K$  for  $\mathbf{x} \in K$ , for all  $K \in \mathcal{T}$ , and  $v_{\mathcal{T}}(\mathbf{x}) = u_\sigma$  for  $\mathbf{x} \in \sigma$ , for all  $\sigma \in \mathcal{E}_{\text{ext}}$ , where  $(u_K)_{K \in \mathcal{T}}, (u_\sigma)_{\sigma \in \mathcal{E}_{\text{ext}}}$  is the solution of (96)-(98) (with (16)-(17)).

### 5.3 Estimate on the solution $(u_{\mathcal{T}}, v_{\mathcal{T}})$

Under the same hypotheses than in Theorem 5.4, let  $(u_{\mathcal{T}}, v_{\mathcal{T}}) \in \mathcal{X}(\mathcal{T})$  be the unique solution of (96)-(98) (see Definition 5.5). We give in this Subsection, a bound for  $(u_{\mathcal{T}}, v_{\mathcal{T}})$ . This will allow us to prove that  $u_{\mathcal{T}}$  converges to  $u \in H^1(\Omega)$  in  $L^2(\Omega)$ -norm.

By multiplying both sides of (96) by  $u_K$  and using the techniques of Section 5.2, we get:

$$(1 + C_{\alpha} C_{32}) |u_{\mathcal{T}}|_{1,\mathcal{T}}^2 + \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} \frac{m(\sigma)}{d_{\sigma}} (1 + d_{\sigma} C_{\alpha}) (u_K - u_{\sigma})^2 + \frac{1}{2} |v_{\mathcal{T}}|_{\alpha, \mathcal{Z}(\mathcal{T})}^2 \leq \mathbb{T}_1^{\mathcal{T}} + \mathbb{T}_2^{\mathcal{T}}, \quad (108)$$

where  $\mathbb{T}_1^{\mathcal{T}}$  and  $\mathbb{T}_2^{\mathcal{T}}$  are defined as in (24). Inequalities (25) and (26) (of the case  $\alpha$  is constant) remain valid here. Thus, the following Theorem holds

**THEOREM 5.6** Assume Assumption 5.1 and let  $C_{\alpha} = \min_{\partial\Omega} \alpha_t$ . Let  $\mathcal{T}$  be an admissible mesh in the sense of Definition 3.1. Then, There exists  $\delta < 0$ , only depending on  $\Omega$ , such that if  $C_{\alpha} > \delta$  the system (96)-(98) (with (16)-(17)), where  $\{(f_K, g_{\sigma}) : (K, \sigma) \in \mathcal{T} \times \mathcal{E}_{\text{ext}}\}$  is given by (9), has a unique solution. Furthermore, one has:

$$|(u_{\mathcal{T}}, v_{\mathcal{T}})|_{1, \mathcal{X}(\mathcal{T})} + |v_{\mathcal{T}}|_{\alpha, \mathcal{Z}(\mathcal{T})} \leq C_{33} \{\|f\|_{0, \Omega} + \|g\|_{0, \partial\Omega}\}, \quad (109)$$

where  $C_{33}$  is only depending on  $\Omega$  and  $\alpha$  (the semi-norms  $|\cdot, \cdot|_{1, \mathcal{X}(\mathcal{T})}$  and  $|\cdot|_{\alpha, \mathcal{Z}(\mathcal{T})}$  are defined in Definition 3.4 and Definition 5.3).

### 5.4 The convergence of $(u_{\mathcal{T}}, v_{\mathcal{T}})$

In this section, we assume Assumption 5.1 and  $C_{\alpha} = \min_{\partial\Omega} \alpha_t > \delta$ , where  $\delta$  is given by Theorem 5.6. For an an admissible mesh  $\mathcal{T}$  in the sense of Definition 3.1, let  $(u_{\mathcal{T}}, v_{\mathcal{T}})$  be the unique solution  $(u_{\mathcal{T}}, v_{\mathcal{T}})$  of (96)-(98) (with (16)-(17)). The objective is to prove the convergence of  $(u_{\mathcal{T}}, v_{\mathcal{T}})$  to the solution of (87)-(88) as the mesh size goes to zero. The proof follows that of the case “ $\alpha$  constant”. We begin with the following Lemma which proof is similar to that of Lemma 4.5.

**LEMMA 5.7** Under the hypotheses of Theorem 5.6, let  $(u_{\mathcal{T}}, v_{\mathcal{T}}) \in \mathcal{X}(\mathcal{T})$  be the solution of (96)-(98) in the sense of Definition 5.5. Then the following estimate holds

$$\|v_{\mathcal{T}}\|_{0, \partial\Omega} \leq C_{34}, \quad (110)$$

where  $C_{34}$  depends only on  $(\Omega, \alpha, f, g)$ .

Since the set  $Y$  of the approximations  $u_{\mathcal{T}}$  is bounded in the  $L^2(\Omega)$  (thanks to the discrete mean Poincaré inequality and Inequality (109)), we are able now to justify that  $u_{\mathcal{T}}$  converges to some  $u$  as  $\text{size}(\mathcal{T})$  goes to zero. Uniform boundedness (109) and compactness result of [EGH 00] in case of Neumann problem yields that the set  $Y$  is relatively compact in  $L^2(\Omega)$ . In addition to this, if a sequence  $u_{\mathcal{T}_n}$  ( $n \in \mathbb{N}$ ) converges to a function  $u$  in  $L^2$ -norm as  $\text{size}(\mathcal{T}_n)$  goes to 0, then  $u \in H^1(\Omega)$ . Furthermore, Lemma 5.7 implies that  $v_{\mathcal{T}_n}$  converges weakly to some  $v \in L^2(\partial\Omega)$ , up to a subsequence. We start by proving:

$$\begin{aligned} - \int_{\Omega} u(\mathbf{x}) \Delta \varphi(\mathbf{x}) d\mathbf{x} &+ \int_{\partial\Omega} \varphi_n(\mathbf{x}) v(\mathbf{x}) d\gamma(\mathbf{x}) = \int_{\Omega} f(\mathbf{x}) \varphi(\mathbf{x}) d\mathbf{x} \\ &+ \int_{\partial\Omega} g(\mathbf{x}) \varphi(\mathbf{x}) d\gamma(\mathbf{x}) + \int_{\partial\Omega} \varphi_t(\mathbf{x}) \alpha(\mathbf{x}) v(\mathbf{x}) d\gamma(\mathbf{x}), \forall \varphi \in \mathcal{C}^2(\overline{\Omega}). \end{aligned} \quad (111)$$

To simplify the notations, we set  $u_{\mathcal{T}_n} = u_{\mathcal{T}}$  and  $v_{\mathcal{T}_n} = v_{\mathcal{T}}$ .

Let  $\varphi \in \mathcal{C}^2(\overline{\Omega})$  and consider the function  $\varphi_{\mathcal{T}} = (\varphi_{\mathcal{T}}^{(1)}, \varphi_{\mathcal{T}}^{(2)}) \in \mathcal{X}(\mathcal{T})$  (see Definition 3.3)

defined by  $\varphi_T^{(1)}(\mathbf{x}) = \varphi_K = \varphi(\mathbf{x}_K)$ , for  $\mathbf{x} \in K$  and for any control volume  $K$ , and  $\varphi_T^{(2)}(\mathbf{x}) = \varphi_\sigma = \varphi(\mathbf{y}_\sigma)$  for  $\mathbf{x} \in \sigma$ , for any  $\sigma \in \mathcal{E}_{\text{ext}}$ . Multiplying both sides of equation (96) by  $\varphi_K$  and using (97) combined with the techniques used for  $\alpha$  constant, we get

$$\begin{aligned} - \int_{\Omega} u_T(\mathbf{x}) \Delta \varphi(\mathbf{x}) d\mathbf{x} &+ \int_{\partial\Omega} \varphi_n(\mathbf{x}) v_T(\mathbf{x}) d\gamma(\mathbf{x}) + \bar{r} = \int_{\Omega} f(\mathbf{x}) \varphi_T^{(1)}(\mathbf{x}) d\mathbf{x} \\ &+ \int_{\partial\Omega} g(\mathbf{x}) \varphi_T^{(2)}(\mathbf{x}) d\gamma(\mathbf{x}) - \sum_{\sigma \in \mathcal{E}_{\text{ext}}} (\alpha(\sigma_e) u_{\sigma,+} - \alpha(\sigma_b) u_{\sigma^-,+}) \varphi_\sigma \\ &- \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} (\alpha(\sigma_e) u_{\sigma,+} - \alpha(\sigma_b) u_{\sigma^-,+}) (\varphi_K - \varphi_\sigma) \\ &+ \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} m(\sigma) (\varphi_K - \varphi_\sigma) g_\sigma, \end{aligned} \quad (112)$$

where  $\bar{r}$ ,  $r(\varphi, T)$  and  $s$  are defined as for  $\alpha$  constant, i.e. as in (29)-(32). We remark that the r.h.s. of (112) differs from the r.h.s. of (31) only by the third and fourth term.

We first prove that the fourth term on the r.h.s. of (112) goes to zero, as  $\text{size}(T) \rightarrow 0$ . Let  $\varphi_{K,\sigma} = \varphi_K - \varphi_\sigma$  and write

$$\sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} (\alpha(\sigma_e) u_{\sigma,+} - \alpha(\sigma_b) u_{\sigma^-,+}) \varphi_{K,\sigma} = \mathbb{T}_7^T + \mathbb{T}_8^T + \mathbb{T}_9^T + \mathbb{T}_{10}^T, \quad (113)$$

where

$$\mathbb{T}_7^T = \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} \alpha(\sigma_e) (u_{\sigma,+} - u_{\sigma,-}) \varphi_{K,\sigma} \text{ and } \mathbb{T}_8^T = \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} \alpha(\sigma_e) u_{\sigma,-} \varphi_{K,\sigma},$$

and

$$\mathbb{T}_9^T = - \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} \alpha(\sigma_b) (u_{\sigma^-,+} - u_{\sigma^-, -}) \varphi_{K,\sigma} \text{ and } \mathbb{T}_{10}^T = - \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} \alpha(\sigma_b) u_{\sigma^-, -} \varphi_{K,\sigma}.$$

To estimate  $\mathbb{T}_7^T$  and  $\mathbb{T}_9^T$ , we assume that the mesh  $T$  satisfies, for some  $(\zeta_3, \zeta_4) \in (\mathbb{R}_*^+)^2$ , the following condition:

$$\zeta_3 m(\sigma) \leq m(\sigma^+) \leq \zeta_4 m(\sigma), \quad \forall \sigma \in \mathcal{E}_{\text{ext}}, \quad (114)$$

where  $\sigma^+$  is defined in Definition 3.2. Thanks to Estimate (109), the fact that  $|\varphi_{K,\sigma}| \leq \text{size}(T) \|\nabla \varphi\|_{L^\infty(\Omega)}$  and under the Assumption that the mesh  $T$  satisfies, for some  $(\zeta_1, \zeta_3, \zeta_4) \in (\mathbb{R}_*^+)^3$ , the conditions (34) and (114), we have the following estimates:

$$|\mathbb{T}_7^T| \leq C_{35} \sqrt{\text{size}(T)} \text{ and } |\mathbb{T}_9^T| \leq C_{36} \sqrt{\text{size}(T)}, \quad (115)$$

where  $C_{35}$  depends on  $(\Omega, \zeta_1, f, g, \varphi, \alpha)$  and  $C_{36}$  depends on  $(\Omega, \zeta_1, \zeta_4, f, g, \varphi, \alpha)$ . It remains to estimate  $\mathbb{T}_8^T + \mathbb{T}_{10}^T$ . Indeed

$$\mathbb{T}_8^T + \mathbb{T}_{10}^T = \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} (\alpha(\sigma_e) - \alpha(\sigma_b)) u_{\sigma,-} \varphi_{K,\sigma} + \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} \alpha(\sigma_b) (u_{\sigma,-} - u_{\sigma^-, -}) \varphi_{K,\sigma}. \quad (116)$$

Thanks to Lemma 5.7 and using  $|\varphi_{K,\sigma}| \leq \text{size}(T) \|\nabla \varphi\|_{L^\infty(\Omega)}$  and  $|\alpha(\sigma_e) - \alpha(\sigma_b)| \leq m((\sigma)) \|\nabla \alpha\|_{L^\infty(\partial\Omega)}$ , the first term on the r.h.s. of (116) can be bounded as:

$$\left| \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} (\alpha(\sigma_e) - \alpha(\sigma_b)) u_{\sigma,-} \varphi_{K,\sigma} \right| \leq C_{37} \text{size}(T), \quad (117)$$

where  $C_{37}$  depends on  $(\Omega, \zeta_3, f, g, \varphi, \alpha)$ . To estimate the second term on r.h.s. of (116), we consider the following sets:

$$\mathbb{Z}_1 = \{\sigma \in \mathcal{E}_{\text{ext}} : \alpha(\sigma_e) \geq 0 \text{ and } \alpha(\sigma_b) \geq 0\}, \quad (118)$$

$$\mathbb{Z}_2 = \{\sigma \in \mathcal{E}_{\text{ext}} : \alpha(\sigma_e) \geq 0 \text{ and } \alpha(\sigma_b) \leq 0\}, \quad (119)$$

$$\mathbb{Z}_3 = \{\sigma \in \mathcal{E}_{\text{ext}} : \alpha(\sigma_e) \leq 0 \text{ and } \alpha(\sigma_b) \leq 0\}. \quad (120)$$

We have:

$$\begin{aligned} \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} \alpha(\sigma_b)(u_{\sigma,-} - u_{\sigma^-, -})\varphi_{K,\sigma} &= \sum_{\mathcal{E}_K \cap \mathbb{Z}_1} (\alpha(\sigma_b) - \alpha(\sigma_e))(u_{\sigma+} - u_{\sigma})\varphi_{K,\sigma} \\ &+ \sum_{\mathcal{E}_K \cap \mathbb{Z}_1} \alpha(\sigma_e)(u_{\sigma+} - u_{\sigma})\varphi_{K,\sigma} \\ &+ \sum_{\mathcal{E}_K \cap \mathbb{Z}_2} \alpha(\sigma_b)(u_{\sigma+} - u_{\sigma})\varphi_{K,\sigma} \\ &+ \sum_{\mathcal{E}_K \cap \mathbb{Z}_2} \alpha(\sigma_b)(u_{\sigma} - u_{\sigma-})\varphi_{K,\sigma} \\ &+ \sum_{\mathcal{E}_K \cap \mathbb{Z}_3} \alpha(\sigma_b)(u_{\sigma} - u_{\sigma-})\varphi_{K,\sigma} \end{aligned} \quad (121)$$

The first term on the r.h.s. of (121) can be bounded using triangular Inequality, the assumption (114) and Lemma 5.7:

$$\left| \sum_{\mathbb{Z}_1} (\alpha(\sigma_b) - \alpha(\sigma_e))(u_{\sigma+} - u_{\sigma})\varphi_{K,\sigma} \right| \leq C_{38} \text{size}(\mathcal{T}), \quad (122)$$

where  $C_{38}$  depends on  $(\alpha, \varphi, \Omega, f, g, \zeta_3)$ . Noting that  $\{u_{\sigma,+}, u_{\sigma,-}\} = \{u_{\sigma+}, u_{\sigma}\}$  and

$$\left| \sum_{\mathbb{Z}_1 \cap \mathcal{E}_K} \alpha(\sigma_e)(u_{\sigma+} - u_{\sigma})\varphi_{K,\sigma} \right| \leq \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} |\alpha(\sigma_e)(u_{\sigma+} - u_{\sigma})\varphi_{K,\sigma}|,$$

the first estimate of (115) can be applied here to get:

$$\left| \sum_{\mathbb{Z}_1 \cap \mathcal{E}_K} \alpha(\sigma_e)(u_{\sigma+} - u_{\sigma})\varphi_{K,\sigma} \right| \leq C_{35} \sqrt{\text{size}(\mathcal{T})}. \quad (123)$$

The techniques used to bound the first and the second term on r.h.s. of (121) can be used to obtain the same bound for the third, fourth and fifth term on the r.h.s. of (121). Thus:

$$\left| \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} \alpha(\sigma_b)(u_{\sigma,-} - u_{\sigma^-, -})\varphi_{K,\sigma} \right| \leq C_{39} \sqrt{\text{size}(\mathcal{T})}, \quad (124)$$

where  $C_{39}$  depends on  $(\Omega, \zeta_3, \zeta_1, \zeta_4, f, g, \varphi, \alpha)$ . This, with (117), (116), implies

$$|\mathbb{T}_8^{\mathcal{T}} + \mathbb{T}_{10}^{\mathcal{T}}| \leq C_{40} \sqrt{\text{size}(\mathcal{T})}, \quad (125)$$

where  $C_{40}$  depends on  $(\Omega, \zeta_3, \zeta_1, \zeta_4, f, g, \varphi, \alpha)$ . Combining now (113), (115) and (125), we get the following estimate for the fourth term on the r.h.s. of (112):

$$\left| \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} (\alpha(\sigma_e)u_{\sigma,+} - \alpha(\sigma_b)u_{\sigma^-, +})\varphi_{K,\sigma} \right| \leq C_{41} \sqrt{\text{size}(\mathcal{T})}, \quad (126)$$

where  $C_{41}$  depends on  $(\Omega, \zeta_3, \zeta_1, \zeta_4, f, g, \varphi, \alpha)$ .



We now turn to estimate the third term on the r.h.s. of (112). Reordering the sum and using formula (10), we get (recall that  $\varphi_\sigma = \varphi(\mathbf{y}_\sigma)$ ):

$$\begin{aligned}
 \sum_{\sigma \in \mathcal{E}_{\text{ext}}} (\alpha(\sigma_e)u_{\sigma,+} - \alpha(\sigma_b)u_{\sigma^-,+})\varphi_\sigma &= \sum_{\sigma \in \mathcal{E}_{\text{ext}}} \alpha(\sigma_e)u_{\sigma,+}(\varphi_\sigma - \varphi_{\sigma^+}) \\
 &= \sum_{\sigma \in \mathcal{E}_{\text{ext}}} \alpha(\sigma_e)u_{\sigma,+}(\varphi(\mathbf{y}_\sigma) - \varphi(\sigma_b)) \\
 &\quad - \sum_{\sigma \in \mathcal{E}_{\text{ext}}} \alpha(\sigma_e)u_{\sigma,+}(\varphi(\mathbf{y}_{\sigma^+}) - \varphi(\sigma_e)) \\
 &\quad - \sum_{\sigma \in \mathcal{E}_{\text{ext}}} \alpha(\sigma_e)u_{\sigma,+}(\varphi(\sigma_e) - \varphi(\sigma_b)) \\
 &= \sum_{\sigma \in \mathcal{E}_{\text{ext}}} \{\varphi(\mathbf{y}_\sigma) - \varphi(\sigma_b)\}(\alpha(\sigma_e)u_{\sigma,+} - \alpha(\sigma_b)u_{\sigma^-,+}) \\
 &\quad - \int_{\partial\Omega} \varphi_t(\mathbf{x})\alpha_T(\mathbf{x})v_T^+(\mathbf{x}) d\gamma(\mathbf{x}), \tag{127}
 \end{aligned}$$

where  $(\alpha_T, v_T^+) \in (\mathcal{Z}(T))^2$  is defined by  $(\alpha_T(\mathbf{x}), v_T^+(\mathbf{x})) = (\alpha(\sigma_e), u_{\sigma,+})$  for  $\mathbf{x} \in \sigma$ , for all  $\sigma \in \mathcal{E}_{\text{ext}}$ . As it is done to obtain the estimate (126), we can obtain the same estimate for the first term on the r.h.s. of (127):

$$\left| \sum_{\sigma \in \mathcal{E}_{\text{ext}}} \{\varphi(\mathbf{y}_\sigma) - \varphi(\sigma_b)\}(\alpha(\sigma_e)u_{\sigma,+} - \alpha(\sigma_b)u_{\sigma^-,+}) \right| \leq C_{42} \sqrt{\text{size}(T)}, \tag{128}$$

where  $C_{42}$  depends on  $(\Omega, \zeta_3, f, g, \zeta_4, \varphi, \alpha)$ . We now turn to the second term on the r.h.s. of (127):

$$\begin{aligned}
 \int_{\partial\Omega} \varphi_t(\mathbf{x})\alpha_T(\mathbf{x})v_T^+(\mathbf{x}) d\gamma(\mathbf{x}) &= \int_{\partial\Omega} \varphi_t(\mathbf{x})\alpha_T(\mathbf{x})(v_T^+ - v_T)(\mathbf{x}) d\gamma(\mathbf{x}) \\
 &\quad + \int_{\partial\Omega} \varphi_t(\mathbf{x})\alpha_T(\mathbf{x})v_T(\mathbf{x}) d\gamma(\mathbf{x}). \tag{129}
 \end{aligned}$$

Since  $\|\alpha_T - \alpha\|_{L^\infty(\Omega)} \rightarrow 0$  and  $\|v_T - v\|_{L^2(\partial\Omega)} \rightarrow 0$ , as  $\text{size}(T) \rightarrow 0$ , then:

$$\int_{\partial\Omega} \varphi_t(\mathbf{x})\alpha_T(\mathbf{x})v_T(\mathbf{x}) d\gamma(\mathbf{x}) \rightarrow \int_{\partial\Omega} \varphi_t(\mathbf{x})\alpha(\mathbf{x})v(\mathbf{x}) d\gamma(\mathbf{x}), \text{ as } \text{size}(T) \rightarrow 0. \tag{130}$$

We now prove that the first term on the r.h.s. of (129) goes to zero. Indeed, using (109) and the fact that  $\{|u_{\sigma,+} - u_\sigma|\} = \{0, |u_{\sigma^+} - u_\sigma|\} = \{0, |u_{\sigma,+} - u_{\sigma,-}|\}$ , we get

$$\left| \int_{\partial\Omega} \varphi_t(\mathbf{x})\alpha_T(\mathbf{x})(v_T^+ - v_T)(\mathbf{x}) d\gamma(\mathbf{x}) \right| \leq C_{43} \sqrt{\text{size}(T)}, \tag{131}$$

where  $C_{43}$  depends on  $(\Omega, f, g, \varphi, \alpha)$ . Let  $\text{size}(T)$  tends to zero in (129). Using (130) and (131), we get:

$$\int_{\partial\Omega} \varphi_t(\mathbf{x})\alpha_T(\mathbf{x})v_T^+(\mathbf{x}) d\gamma(\mathbf{x}) \rightarrow \int_{\partial\Omega} \varphi_t(\mathbf{x})\alpha(\mathbf{x})v(\mathbf{x}) d\gamma(\mathbf{x}), \text{ as } \text{size}(T) \rightarrow 0. \tag{132}$$

This, with (127) and (128), implies

$$\sum_{\sigma \in \mathcal{E}_{\text{ext}}} (\alpha(\sigma_e)u_{\sigma,+} - \alpha(\sigma_b)u_{\sigma^-,+})\varphi_\sigma \rightarrow - \int_{\partial\Omega} \varphi_t(\mathbf{x})\alpha(\mathbf{x})v(\mathbf{x}) d\gamma(\mathbf{x}), \text{ as } \text{size}(T) \rightarrow 0. \tag{133}$$

Writing (112) with  $\mathcal{T} = \mathcal{T}_n$ , passing to the limit as  $n$  tends to infinity (we assume that  $\text{size}(\mathcal{T}_n) \rightarrow 0$ , as  $n \rightarrow \infty$ ) and using (37), (126) and (133), we get equation (111). Using Lemma 4.6, we get  $\tilde{\gamma}(u) = v$  a.e. on  $\partial\Omega$ . This with an integration by part in (111) implies that, for any  $\varphi \in \mathcal{D}(\overline{\Omega})$ , we have:

$$\begin{aligned} \int_{\Omega} \nabla u(\mathbf{x}) \cdot \nabla \varphi(\mathbf{x}) d\mathbf{x} &= \int_{\Omega} f(\mathbf{x}) \varphi(\mathbf{x}) d\mathbf{x} + \int_{\partial\Omega} g(\mathbf{x}) \varphi(\mathbf{x}) d\gamma(\mathbf{x}) \\ &+ \int_{\partial\Omega} \varphi_t(\mathbf{x}) \alpha(\mathbf{x}) \tilde{\gamma}(u)(\mathbf{x}) d\gamma(\mathbf{x}), \forall \varphi \in \mathcal{C}^2(\overline{\Omega}). \end{aligned}$$

This, with formula (5), gives:

$$\begin{aligned} \int_{\Omega} \nabla u(\mathbf{x}) \cdot \nabla \varphi(\mathbf{x}) d\mathbf{x} &+ \int_{\Omega} (\varphi_x(\mathbf{x}) (\alpha u)_y(\mathbf{x}) - \varphi_y(\mathbf{x}) (\alpha u)_x(\mathbf{x})) d\mathbf{x} \\ &= \int_{\Omega} f(\mathbf{x}) \varphi(\mathbf{x}) d\mathbf{x} + \int_{\partial\Omega} g(\mathbf{x}) \varphi(\mathbf{x}) d\gamma(\mathbf{x}), \forall \varphi \in \mathcal{D}(\overline{\Omega}). \end{aligned} \quad (134)$$

Thanks to the density  $\mathcal{D}(\overline{\Omega})$  in  $H^1(\Omega)$ , the formulation (134) is equivalent to (88). Letting  $\text{size}(\mathcal{T})$  tends to zero in (98) and using the fact that  $u_{\mathcal{T}}$  tends to  $u$  in the  $L^2(\Omega)$ -norm, we get (87). Since the solution  $u$  of (87)-(88) is unique, the whole family  $u_{\mathcal{T}}$  converges to the solution  $u \in H^1(\Omega)$  of (87)-(88) in  $L^2(\Omega)$  and the whole family  $v_{\mathcal{T}}$  converges to  $\tilde{\gamma}(u)$  for the weak topology of  $L^2(\partial\Omega)$ , as  $\text{size}(\mathcal{T})$  goes to 0.

We now obtain a similar result to that of (45). Multiplying both sides of equation (96) by  $u_K$ ,  $K \in \mathcal{T}$ , summing over  $K \in \mathcal{T}$ , using equation (97) and techniques of (100)-(102) yields

$$\|(u_{\mathcal{T}}, v_{\mathcal{T}})\|_{\star}^2 = \sum_{K \in \mathcal{T}} m(K) f_K u_K + \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} m(\sigma) g_{\sigma} u_K, \quad (135)$$

where  $\|(\cdot, \cdot)\|_{\star}$  is the semi-norm defined by

$$\|(u_{\mathcal{T}}, v_{\mathcal{T}})\|_{\star}^2 = |(u_{\mathcal{T}}, v_{\mathcal{T}})|_{1, \mathcal{X}(\mathcal{T})}^2 + \frac{1}{2} |v_{\mathcal{T}}|_{\alpha, \mathcal{Z}(\mathcal{T})}^2 + \frac{1}{2} \int_{\partial\Omega} \alpha_t(\mathbf{x}) v_{\mathcal{T}}^2(\mathbf{x}) d\gamma(\mathbf{x}), \quad (136)$$

where  $|\cdot|_{1, \mathcal{X}(\mathcal{T})}^2$  (resp.  $|\cdot|_{\alpha, \mathcal{Z}(\mathcal{T})}^2$ ) is defined in Definition 3.4 (resp. 5.3)

*Remark 10* We saw in Section 5.2 that the r.h.s. of (136) is nonnegative, thanks to the condition  $C_{\alpha} > \delta$  given in Theorem 5.6 (this condition is assumed in all the present section).

On the other hand, if we replace  $v$  with  $u$  in (88), we get:

$$|u|_{1, \Omega}^2 + \frac{1}{2} \int_{\partial\Omega} \alpha_t(\mathbf{x}) u^2(\mathbf{x}) d\gamma(\mathbf{x}) = \int_{\Omega} f(\mathbf{x}) u(\mathbf{x}) d\mathbf{x} + \int_{\partial\Omega} g(\mathbf{x}) \tilde{\gamma}(u)(\mathbf{x}) d\gamma(\mathbf{x}). \quad (137)$$

Letting  $\text{size}(\mathcal{T})$  tends to zero in the r.h.s. of (135), we get:

$$\|(u_{\mathcal{T}}, v_{\mathcal{T}})\|_{\star}^2 \rightarrow \int_{\Omega} f(\mathbf{x}) u(\mathbf{x}) d\mathbf{x} + \int_{\partial\Omega} g(\mathbf{x}) \tilde{\gamma}(u)(\mathbf{x}) d\gamma(\mathbf{x}), \quad (138)$$

this, with (137), implies:

$$\|(u_{\mathcal{T}}, v_{\mathcal{T}})\|_{\star}^2 \rightarrow |u|_{1, \Omega}^2 + \frac{1}{2} \int_{\partial\Omega} \alpha_t(\mathbf{x}) u^2(\mathbf{x}) d\gamma(\mathbf{x}) \text{ as } \text{size}(\mathcal{T}) \rightarrow 0. \quad (139)$$

Up to now, we obtained the convergence of the approximate solution  $(u_{\mathcal{T}}, v_{\mathcal{T}})$  when the solution  $u$  of (87)-(88) only satisfies  $u \in H^1(\Omega)$ . Assume now that the weak solution  $u$  of (87)-(88) satisfies  $u \in \mathcal{C}^2(\overline{\Omega})$ . Using the same techniques as in the Sections 4.4, 5.2 and 5.3

(with the same notations), there exists  $\delta < 0$  only depending on  $\Omega$  ( $\delta$  is as in Theorem 5.6) and there exists  $C_{44}$ , only depending on  $(\alpha, \Omega)$ , such that for  $C_\alpha > \delta$  we have:

$$\begin{aligned} |(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}})|_{1, \mathcal{X}(\mathcal{T})}^2 + |\bar{e}_{\mathcal{T}}|_{\alpha, \mathcal{Z}(\mathcal{T})}^2 &\leq C_{44} \{ |\sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} m(\sigma) R_{K, \sigma} e_K| + |\sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} m(\sigma) R_{K, \sigma} e_\sigma| \\ &+ |\sum_{\sigma \in \mathcal{E}_{\text{ext}}} (r_\sigma - r_{\sigma^-}) e_\sigma| \}, \end{aligned} \quad (140)$$

where

$$r_\sigma = \alpha(\sigma_e) u(\sigma_e) - \alpha(\sigma_e) u(\mathbf{y}_{\sigma, +}), \quad (141)$$

and  $R_{K, \sigma}$  is defined by (53)-(54) and  $(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}}) \in \mathcal{X}(\mathcal{T})$  is defined by  $(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}}) = (u(\mathbf{x}_K) - u_K, u(\mathbf{y}_{\sigma, +}) - u_{\sigma, +})$  (recall that  $\mathbf{y}_{\sigma, +} = \mathbf{y}_\sigma$  (resp.  $\mathbf{y}_{\sigma, +}$ ) if  $\alpha(\sigma_e) \geq 0$  (resp.  $\alpha(\sigma_e) < 0$ )). Since  $u \in \mathcal{C}^2(\bar{\Omega})$ , then:

$$|r_\sigma| \leq C_{45} \text{size}(\mathcal{T}), \quad (142)$$

where  $C_{45}$  depends on  $(\alpha, u)$ . If, we assume that the condition (72) is fulfilled for some positive number  $\zeta_2$  (this implies that the conditions (34) and (114) are fulfilled). Then, Inequality (140), with (55) and (142), implies:

$$|(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}})|_{1, \mathcal{X}(\mathcal{T})}^2 + |\bar{e}_{\mathcal{T}}|_{\alpha, \mathcal{Z}(\mathcal{T})}^2 \leq C_{46} \text{size}(\mathcal{T}), \quad (143)$$

where  $C_{46}$  depends on  $(\alpha, \Omega, \zeta_2, u)$ . Now, we use the techniques of Section 4.4 and Error Estimate (143) to obtain the following error estimate in  $L^2(\Omega)$ -norm:

$$\|e_{\mathcal{T}}\|_{L^2(\Omega)}^2 \leq C_{47} \text{size}(\mathcal{T}), \quad (144)$$

where  $C_{47}$  depends on  $(\alpha, \Omega, \zeta_2, u)$ . The following Theorem summarizes the results of the oblique derivative problem with a sufficiently smooth function  $\alpha$ .

**THEOREM 5.8 (CONVERGENCE AND ERROR ESTIMATE FOR SMOOTH  $\alpha$ )** Assume Assumption 5.1 and let  $C_\alpha = \min_{\partial\Omega} \alpha_t$ . Let  $\mathcal{T}$  be an admissible mesh in the sense of Definition 3.1 and let  $\{(f_K, g_\sigma) : (K, \sigma) \in \mathcal{T} \times \mathcal{E}_{\text{ext}}\}$  be given by (9). Then, There exists  $\delta < 0$ , only depending on  $\Omega$ , such that if  $C_\alpha > \delta$  one has:

1. Problem (87)-(88) has a unique solution  $u \in H^1(\Omega)$ ,
2. System (96)-(98) has a unique solution  $(u_{\mathcal{T}}, v_{\mathcal{T}}) \in \mathcal{X}(\mathcal{T})$  in the sense of Definition 5.5.

If in addition, the conditions (114) and (34) are fulfilled for some  $(\zeta_1, \zeta_3, \zeta_4) \in (\mathbb{R}_+^*)^3$ . Then:

$$u_{\mathcal{T}} \rightarrow u \text{ in } L^2(\Omega), \text{ as } \text{size}(\mathcal{T}) \rightarrow 0, \quad (145)$$

$$\|(u_{\mathcal{T}}, v_{\mathcal{T}})\|_\star^2 \rightarrow |u|_{1, \Omega}^2 + \frac{1}{2} \int_{\partial\Omega} \alpha_t(\mathbf{x}) u^2(\mathbf{x}) d\gamma(\mathbf{x}), \text{ as } \text{size}(\mathcal{T}) \rightarrow 0, \quad (146)$$

and

$$v_{\mathcal{T}} \rightarrow \tilde{\gamma}(u) \text{ in } L^2(\partial\Omega) \text{ for the weak topology, as } \text{size}(\mathcal{T}) \rightarrow 0, \quad (147)$$

where the semi-norm  $\|(\cdot, \cdot)\|_\star$  is defined in (136) and  $\tilde{\gamma}$  is the classical trace operator from  $H^1(\Omega)$  to  $L^2(\partial\Omega)$ . Assume furthermore that:

3. The weak solution of (87)-(88) satisfies  $u \in \mathcal{C}^2(\bar{\Omega})$ .
4. The condition (72) is fulfilled for some positive number  $\zeta_2$  (this implies that the conditions (34) and (114) are fulfilled).

Then the following error estimates hold:

$$\|e_{\mathcal{T}}\|_{L^2(\Omega)}^2 \leq C_{47} \text{size}(\mathcal{T}), \quad (148)$$

$$|(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}})|_{1, \mathcal{X}(\mathcal{T})}^2 + |\bar{e}_{\mathcal{T}}|_{\alpha, \mathcal{Z}(\mathcal{T})}^2 \leq C_{46} \text{size}(\mathcal{T}), \quad (149)$$

where  $(C_{46}, C_{47})$  depend on  $(\alpha, \Omega, \zeta_2, u)$ ; and  $(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}}) \in \mathcal{X}(\mathcal{T})$  is defined by  $(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}}) = (u(\mathbf{x}_K) - u_K, u(\mathbf{y}_{\sigma,+}) - u_{\sigma,+})$ , on  $(K, \sigma) \in \mathcal{T} \times \mathcal{E}_{\text{ext}}$ .

Under the hypotheses of Theorem 5.8, Error Estimate (149) yields the following error estimate:

$$\begin{aligned} & \sum_{\substack{\sigma \in \mathcal{E}_{\text{int}} \\ \sigma = K|L}} m(\sigma) d_{\sigma} \left( \frac{u_L - u_K}{d_{\sigma}} - \frac{1}{m(\sigma)} \int_{\sigma} \nabla u(\mathbf{x}) \cdot \mathbf{n}_{K,\sigma} d\gamma(\mathbf{x}) \right)^2 \\ & + \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} m(\sigma) d_{\sigma} \left( \frac{u_{\sigma} - u_K}{d_{\sigma}} - \frac{1}{m(\sigma)} \int_{\sigma} \nabla u(\mathbf{x}) \cdot \mathbf{n}_{K,\sigma} d\gamma(\mathbf{x}) \right)^2 \leq C_{46} \text{size}(\mathcal{T}). \end{aligned} \quad (150)$$

## 6 Error estimate when $\alpha$ is piecewise constant

In this Section, we consider the case where  $\alpha$  is constant on each line of the boundary  $\partial\Omega$ . The boundary oblique derivative problem we want to present arises, in some cases, from the Laplace equation with Dirichlet boundary conditions (see Remark 12 and [B 05]). We will assume the existence of “sufficiently smooth” solution, that is a function which satisfies (151) and the hypotheses of Assumption 6.1 (equation (151) and Assumption 6.1 will be given below). Then, the uniqueness of such a solution will be deduced from the fact that it is the unique limit of approximate solutions given by a finite volume scheme (see Theorem 6.6). This uniqueness can also be proved, as performed in the second item of Remark 13, by computing the integral  $\int_{\Omega} \Delta u u d\mathbf{x}$ , where  $u$  is a “smooth” solution. The existence of a weak solution  $u \in H^1(\Omega)$  along with a weak formulation for (151), is given in [G 85, Lemma 4.4.4.2]. In some particular cases on  $\alpha$  and under the assumption that  $u \in W_p^1(\Omega)$  with  $p > 1$ , uniqueness results for  $u$  are also presented in [G 85], see for instance [G 85, Lemma 4.4.4.3].

### 6.1 The problem to be solved

We consider the following case of Problem (1):

$$\begin{cases} -\Delta u(\mathbf{x}) = f(\mathbf{x}), \text{ on } \Omega \\ u_n(\mathbf{x}) + (\alpha u)_t(\mathbf{x}) = g(\mathbf{x}) \text{ on } \Gamma = \partial\Omega, \end{cases} \quad (151)$$

where  $\alpha$  is constant on each “line” of the boundary  $\partial\Omega$ . To be more precise, we call  $\{\bar{\Gamma}_j, j = 1, \dots, N\}$  (recall that  $\bar{\Gamma}_j$  denotes the closure of  $\Gamma_j$ ) the lines constituting  $\partial\Omega$ . For each  $j \in \{1, \dots, N\}$ , one has  $\Gamma_j = (\mathcal{S}_{j-1}, \mathcal{S}_j)$ , where  $|\mathcal{S}_j - \mathcal{S}_{j-1}| \mathbf{t}_j = \mathcal{S}_j - \mathcal{S}_{j-1}$ ,  $\mathbf{t}_j = (-\mathbf{n}_j)_y, (\mathbf{n}_j)_x)^t$  and  $\mathbf{n}_j = ((\mathbf{n}_j)_x, (\mathbf{n}_j)_y)^t$  is the normal vector to  $\Gamma_j$ , outward  $\Omega$  (this is similar to the definition of  $\sigma_e$  and  $\sigma_b$  in Definition 3.2). One also has  $\mathcal{S}_N = \mathcal{S}_0$ . The function  $\alpha$  is constant on each  $\Gamma_j$ . Let  $\alpha_j$  be the value of  $\alpha$  on  $\Gamma_j$ .

In order to get an error estimate, we need the following Assumption:

**ASSUMPTION 6.1** We assume that there exists a function  $u \in \mathcal{C}^2(\bar{\Omega})$  satisfying the first equation of (151) for all  $\mathbf{x} \in \Omega$ , the second equation of (151) for all  $\mathbf{x} \in \Gamma_j$  and  $j \in \{1, \dots, N\}$  and such that:

1.  $\int_{\Omega} u(\mathbf{x}) d\mathbf{x} = 0$ ,
2.  $u$  takes the same value on the corners of  $\Omega$ , i.e.  $u(\mathcal{S}_{j-1}) = u(\mathcal{S}_j)$  for all  $j \in \{1, \dots, N\}$ , where  $\mathcal{S}_0 = \mathcal{S}_N$ .

*Remark 11* A consequence of Assumption 6.1 is  $f \in \mathcal{C}(\Omega)$ ,  $g \in \mathcal{C}^1(\Gamma_j)$ , for all  $j \in \{1, \dots, N\}$ , and

$$\int_{\Omega} f(\mathbf{x}) d\mathbf{x} + \int_{\partial\Omega} g(\mathbf{x}) d\gamma(\mathbf{x}) = 0. \quad (152)$$

*Remark 12* Such problems which satisfy Assumption 6.1 arise, for example, when  $u = v_x$  and  $v$  is the solution of the following Dirichlet problem:

$$\begin{cases} -\Delta v(\mathbf{x}) = \bar{f}(\mathbf{x}), \text{ on } \Omega \\ v(\mathbf{x}) = 0, \mathbf{x} \in \partial\Omega, \end{cases} \quad (153)$$

We can see that the derivative  $v_x$  w.r.t.  $x$  of the solution  $v$  of (153), provided that  $v \in \mathcal{C}^3(\overline{\Omega})$ , satisfies an equation like (151) in which the data  $(f, g, \alpha)$  are defined by:

$$(f, g, \alpha) = (\bar{f}_x, -\mathbf{n}_x f, \frac{\mathbf{n}_y}{\mathbf{n}_x}) \quad (154)$$

where  $(-\mathbf{n}_y, \mathbf{n}_x)^t$  are the components of the tangential vector  $\mathbf{t}$  (recall that the components of the normal vector  $\mathbf{n}$ , outward  $\Omega$ , are  $(\mathbf{n}_x, \mathbf{n}_y)^t$ ). We can see that  $u = v_x$  satisfies Assumption 6.1.

*Remark 13* (Assumption 6.1 and uniqueness)

1. In the case of problem (151) with only the first item of Assumption 6.1 (without second item of Assumption 6.1), we have no uniqueness, in general, as it is shown by the following particular case of problem (151):  
Let  $\Omega = (0, 1)^2$ ,  $(f, g) = (0, 0)$  and  $\alpha$  be given by:

$$\alpha = \begin{cases} -1, & \text{on } (0, 1) \times \{0\}, \\ 1, & \text{on } \{1\} \times (0, 1), \\ -1, & \text{on } (0, 1) \times \{1\}, \\ 1, & \text{on } \{0\} \times (0, 1). \end{cases} \quad (155)$$

We can see that the functions  $\{x - y, 0\}$  are two solutions for the problem (151) and they also satisfy the first item of Assumption 6.1 (note that the function  $x - y$  does not satisfy the second item of Assumption 6.1).

2. The uniqueness of a solution which satisfies problem (151) and Assumption 6.1 can also be proved by computing the integral  $\int_{\Omega} \Delta u u d\mathbf{x}$ , see for instance, proof of [G 85, Lemma 4.4.4.3]. Let  $u$  be a solution for the problem (151) with  $(f, g) = (0, 0)$ . Multiplying both sides of the first equation of (151) by  $u$ , using an integration by parts, and using second equation of (151) and (10), we get:

$$\int_{\Omega} |\nabla u|^2 d\mathbf{x} - \frac{1}{2} \sum_{j=1}^N \alpha_j (u(\mathcal{S}_{j+1}) - u(\mathcal{S}_j))^2 = 0, \quad (156)$$

this with Assumption 6.1 implies that  $u = 0$ .

## 6.2 The finite volume scheme for (151)

*Remark 14* When there exists  $j_0 \in \{1, \dots, N\}$  such that  $\alpha_{j_0} = 0$ , the finite volume scheme on the line  $\Gamma_{j_0}$  can be handled as in Neumann problem in [EGH 00].

We assume, for the sake of simplicity, that

$$\alpha_j \neq 0, \forall j \in \{1, \dots, N\}. \quad (157)$$

Therefore, the following sets are considered

$$\mathcal{S}^+ = \{\Gamma_j : \alpha_j > 0\} \text{ and } \mathcal{S}^- = \{\Gamma_j : \alpha_j < 0\}. \quad (158)$$

We need the following Definition:

**DEFINITION 6.2** (Definition of the beginning and the end mesh points on each line  $\Gamma_j$ ) Let  $\mathcal{T}$  be an admissible mesh in the sense of Definition 3.1. For each  $j \in \{1, \dots, N\}$ , we define the beginning mesh point (resp. end mesh point) of the line  $\Gamma_j$ , and we denote it by  $\mathbf{y}_{\text{beg}}^j$  (resp.  $\mathbf{y}_{\text{end}}^j$ ), the point  $\mathbf{y}_\sigma$  such that  $\mathbf{y}_\sigma \in \Gamma_j$  and  $\mathbf{y}_{\sigma-} \notin \Gamma_j$  (resp.  $\mathbf{y}_\sigma \in \Gamma_j$  and  $\mathbf{y}_{\sigma+} \notin \Gamma_j$ ) (note that  $\Gamma_j$  is oriented in the positive direction, see Definition 3.2). We denote then by  $\sigma_{\text{beg}}^j$  (resp.  $\sigma_{\text{end}}^j$ ) the edge  $\sigma \in \mathcal{E}_{\text{ext}}$  which satisfies  $\mathbf{y}_{\text{beg}}^j \in \sigma$  (resp.  $\mathbf{y}_{\text{end}}^j \in \sigma$ ).

To analyze the convergence of the finite volume approximation, we need to use the semi-norm of Definition 3.4 and the following semi-norm:

**DEFINITION 6.3** (A semi-norm on  $\mathcal{Z}(\mathcal{T})$ ) Let  $v_{\mathcal{T}} \in \mathcal{Z}(\mathcal{T})$  (see Definition 3.3) and let  $u_\sigma$  be the value of  $v_{\mathcal{T}}$  on  $\sigma$ , for all  $\sigma \in \mathcal{E}_{\text{ext}}$ . We define the following semi-norm on  $\mathcal{Z}(\mathcal{T})$ :

$$\begin{aligned} |v_{\mathcal{T}}|_{(\alpha_j), \mathcal{Z}(\mathcal{T})}^2 &= \frac{1}{2} \sum_{\Gamma_j \in \mathcal{S}^+} \alpha_j \left( (u_{\text{beg}}^j - u_{\text{end}}^j)^2 + \sum_{(\mathbf{y}_\sigma, \mathbf{y}_{\sigma-}) \in (\Gamma_j)^2} (u_\sigma - u_{\sigma-})^2 \right) \\ &- \frac{1}{2} \sum_{\Gamma_j \in \mathcal{S}^-} \alpha_j \left( (u_{\text{end}}^j - u_{\text{beg}}^j)^2 + \sum_{(\mathbf{y}_\sigma, \mathbf{y}_{\sigma+}) \in (\Gamma_j)^2} (u_{\sigma+} - u_\sigma)^2 \right), \end{aligned} \quad (159)$$

where  $u_{\text{beg}}^j$  (resp.  $u_{\text{end}}^j$ ) denotes the value taken by  $v_{\mathcal{T}}$  on  $\sigma_{\text{beg}}^j$  (resp.  $\sigma_{\text{end}}^j$ ) (see Definition 6.2) (recall that the notations  $u_{\sigma-}$  and  $u_{\sigma+}$  are defined in Definition 3.2).

Let  $\mathcal{T}$  be an admissible mesh in the sense of Definition 3.1 and let  $(u_K)_{K \in \mathcal{T}}, (u_\sigma)_{\sigma \in \mathcal{E}_{\text{ext}}}$  denote the discrete unknowns. The numerical scheme is defined by the following set of equations (recall that  $\{(f_K, g_\sigma), (K, \sigma) \in \mathcal{T} \times \mathcal{E}_{\text{ext}}\}$  is given by (9)):

$$\sum_{\mathcal{E}_K} F_{K,\sigma} = m(K) f_K, \forall K \in \mathcal{T}, \quad (160)$$

where  $F_{K,\sigma}$  is defined by (16)-(17), and, on each  $\Gamma_j \in \mathcal{S}^+$ ,

$$\tau_\sigma(u_\sigma - u_K) = -\alpha_j(u_\sigma - u_{\sigma-}) + m(\sigma)g_\sigma, \forall \sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}} \text{ such that } (\mathbf{y}_\sigma, \mathbf{y}_{\sigma-}) \in (\Gamma_j)^2, \quad (161)$$

$$\tau_\sigma(u_{\text{beg}}^j - u_K) = -\alpha_j(u_{\text{beg}}^j - u_{\text{end}}^j) + m(\sigma_{\text{beg}}^j)g_{\sigma_{\text{beg}}^j}, \quad (162)$$

and, for each  $\Gamma_j \in \mathcal{S}^-$ ,

$$\tau_\sigma(u_\sigma - u_K) = -\alpha_j(u_{\sigma+} - u_\sigma) + m(\sigma)g_\sigma, \forall \sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}} \text{ such that } (\mathbf{y}_\sigma, \mathbf{y}_{\sigma+}) \in (\Gamma_j)^2, \quad (163)$$

$$\tau_\sigma(u_{\text{end}}^j - u_K) = -\alpha_j(u_{\text{beg}}^j - u_{\text{end}}^j) + m(\sigma_{\text{end}}^j)g_{\sigma_{\text{end}}^j}. \quad (164)$$

The equation  $\int_{\Omega} u(\mathbf{x}) d\mathbf{x} = 0$  can be discretized in the following way:

$$\sum_{K \in \mathcal{T}} m(K)u_K = 0. \quad (165)$$

*Remark 15* (Discrete compatibility condition) Summing equation (160) over  $K \in \mathcal{T}$ , we get:

$$\sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma} = \sum_{K \in \mathcal{T}} m(K)f_K, \quad (166)$$

this, with equations (161)-(164), implies:

$$- \sum_{\sigma \in \mathcal{E}_{\text{ext}}} m(\sigma)g_\sigma = \sum_{K \in \mathcal{T}} m(K)f_K. \quad (167)$$

Thus, if there exists a finite volume solution for (160)-(164), then  $\{f_K, K \in \mathcal{T}\}$  and  $\{g_\sigma, \sigma \in \mathcal{E}_{\text{ext}}\}$  should be related by the following discrete compatibility condition:

$$\sum_{K \in \mathcal{T}} m(K)f_K + \sum_{\sigma \in \mathcal{E}_{\text{ext}}} m(\sigma)g_\sigma = 0 \quad (168)$$

Note that this condition is ensured thanks to the compatibility condition (152) (see Remark 11).

### 6.3 Existence and uniqueness of the finite volume solution

To prove the existence and uniqueness of the solution of (160)-(165), we need the following Lemma:

LEMMA 6.4 The following equalities hold, for any  $j \in \{1, \dots, N\}$

$$\begin{aligned} (u_{\text{beg}}^j - u_{\text{end}}^j)u_{\text{beg}}^j + \sum_{(\mathbf{y}_\sigma, \mathbf{y}_{\sigma-}) \in (\Gamma_j)^2} (u_\sigma - u_{\sigma-})u_\sigma &= \frac{1}{2}(u_{\text{beg}}^j - u_{\text{end}}^j)^2 \\ &+ \frac{1}{2} \sum_{(\mathbf{y}_\sigma, \mathbf{y}_{\sigma-}) \in (\Gamma_j)^2} (u_\sigma - u_{\sigma-})^2, \end{aligned} \quad (169)$$

and

$$\begin{aligned} (u_{\text{beg}}^j - u_{\text{end}}^j)u_{\text{end}}^j + \sum_{(\mathbf{y}_\sigma, \mathbf{y}_{\sigma+}) \in (\Gamma_j)^2} (u_{\sigma+} - u_\sigma)u_\sigma &= -\frac{1}{2}(u_{\text{beg}}^j - u_{\text{end}}^j)^2 \\ &- \frac{1}{2} \sum_{(\mathbf{y}_\sigma, \mathbf{y}_{\sigma+}) \in (\Gamma_j)^2} (u_{\sigma+} - u_\sigma)^2. \end{aligned} \quad (170)$$

Let us justify the existence and uniqueness of the solution of (160)-(165). The Proof is mainly the same one of Theorem 4.4.

Let  $M_1$  be the number of elements of  $\mathcal{T}$ ,  $M_2$  the number of elements of  $\mathcal{E}_{\text{ext}}$  and  $M = M_1 + M_2$ . The system (160)-(165) can be viewed as a system of  $M$  unknowns (which are  $\{u_K, K \in \mathcal{T}\}$  and  $\{u_\sigma, \sigma \in \mathcal{E}_{\text{ext}}\}$ ) with  $M$  equations. After the choice of an order for the unknowns and the equations, it can be written as  $Aw = b$ , where  $A$  is  $M \times M$  matrix,

$w \in \mathbb{R}^M$  is the unknown vector and  $b \in \mathbb{R}^M$  is given by the data (namely  $f$  and  $g$ ). Assume that  $b = 0$  (that is  $f_K = 0$  for all  $K \in \mathcal{T}$  and  $g_\sigma = 0$  for all  $\sigma \in \mathcal{E}_{\text{ext}}$ ). Multiplying both sides of equation (160) by  $u_K$ , using the fact that  $\{\mathcal{S}^+, \mathcal{S}^-\}$  is a partition of  $\partial\Omega$  and using (161)-(164) and Lemma 6.4, we get:

$$|(u_{\mathcal{T}}, v_{\mathcal{T}})|_{1, \mathcal{X}(\mathcal{T})}^2 + |v_{\mathcal{T}}|_{(\alpha_j), \mathcal{Z}(\mathcal{T})}^2 = 0, \quad (171)$$

where  $|\cdot|_{(\alpha_j), \mathcal{Z}(\mathcal{T})}$  is defined in Definition 6.3.

Following Remark 1, one deduces that there exists  $C \in \mathbb{R}$  such that  $u_K = C$  for all  $K \in \mathcal{T}$  and  $u_\sigma = C$  for all  $\sigma \in \mathcal{E}_{\text{ext}}$ . This proves that the dimension of the null space of  $A$  is 1. Therefore, the dimension of the range of  $A$  is  $M - 1$ . Since the equality (168) (see Remark 15) is a necessary condition for (160)-(165) to have a solution, it is also a sufficient condition. Furthermore, under this condition on  $f$  and  $g$  (which is given by (152), see Remarks 11 and 15), since the null space of  $A$  is reduced to the set of constant vectors, the system (160)-(165) has a unique solution.

This well posedness of the algebraic system (160)-(165) allows us to introduce the following Definition which is similar to that when  $\alpha$  is constant or a smooth function.

**DEFINITION 6.5** An element  $(u_{\mathcal{T}}, v_{\mathcal{T}}) \in \mathcal{X}(\mathcal{T})$  (see Definition 3.3) is a solution of (160)-(165) if  $u_{\mathcal{T}}(\mathbf{x}) = u_K$  for  $\mathbf{x} \in K$ , for all  $K \in \mathcal{T}$ , and  $v_{\mathcal{T}}(\mathbf{x}) = u_\sigma$  for  $\mathbf{x} \in \sigma$ , for all  $\sigma \in \mathcal{E}_{\text{ext}}$ , where  $(u_K)_{K \in \mathcal{T}}, (u_\sigma)_{\sigma \in \mathcal{E}_{\text{ext}}}$  is the solution of (160)-(165).

## 6.4 Error estimate

In this Section, we show, under the Assumption 6.1, that  $u_{\mathcal{T}}$  converges to the solution  $u$  of (151) given by Assumption 6.1. Furthermore, the order of convergence is  $\sqrt{\text{size}(\mathcal{T})}$ . To prove this result, we follow the same steps as in Section 4.4. Let  $C_{\mathcal{T}} \in \mathbb{R}$  be such that

$$\sum_{K \in \mathcal{T}} m(K) \bar{u}(\mathbf{x}_K) = 0, \quad (172)$$

where  $\bar{u} = u + C_{\mathcal{T}}$ .

For each  $(K, \sigma) \in \mathcal{T} \times \mathcal{E}_{\text{ext}}$ , let  $e_K = \bar{u}(\mathbf{x}_K) - u_K$  and  $e_\sigma = \bar{u}(\mathbf{y}_\sigma) - u_\sigma$ , where  $u_K$  (resp.  $u_\sigma$ ) is the value of  $u_{\mathcal{T}}$  (resp.  $v_{\mathcal{T}}$ ) on  $K$  (resp.  $\sigma$ ) (recall that  $\mathbf{y}_\sigma$  is defined in Definition 3.1).

We consider  $(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}}) \in \mathcal{X}(\mathcal{T})$  defined by  $(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}}) = (e_K, e_\sigma)$ , on  $(K, \sigma) \in \mathcal{T} \times \mathcal{E}_{\text{ext}}$ .

One defines  $R_{K, \sigma}$  as in (53)-(54) and then (thanks to Assumption 6.1, first item) the estimate (55) holds.

Since  $-\Delta \bar{u} = f$ , then (56)-(59) remain valid here.

Using (58) to get, on  $\Gamma_j \in \mathcal{S}^+$  and for  $\sigma \in \mathcal{E}_{\text{ext}}$  such that  $(\mathbf{y}_\sigma, \mathbf{y}_{\sigma^-}) \in (\Gamma_j)^2$ , we have (note that, we will denote  $\sigma = (\sigma_b, \sigma_e)$  in the positive orientation, see Definition 3.2)

$$\tau_\sigma(\bar{u}(\mathbf{y}_\sigma) - \bar{u}(\mathbf{x}_K)) + \alpha_j(\bar{u}(\mathbf{y}_\sigma) - \bar{u}(\mathbf{y}_{\sigma^-})) + r_{\sigma^-}^{j,+} - r_{\sigma^+}^{j,+} = m(\sigma)g_\sigma + m(\sigma)R_{K, \sigma}, \quad (173)$$

where

$$r_{\sigma^-}^{j,+} = \alpha_j(\bar{u}(\sigma_e) - \bar{u}(\mathbf{y}_\sigma)), \quad (174)$$

and (the case  $\mathbf{y}_\sigma \in \Gamma_j$  and  $\mathbf{y}_{\sigma^-} \notin \Gamma_j$ )

$$\begin{aligned} \tau_\sigma(\bar{u}(\mathbf{y}_{\text{beg}}^j) - \bar{u}(\mathbf{x}_K)) + \alpha_j(\bar{u}(\mathbf{y}_{\text{beg}}^j) - \bar{u}(\mathbf{y}_{\text{end}}^j)) + \alpha_j(\bar{u}(\sigma_e) - \bar{u}(\mathbf{y}_{\text{beg}}^j)) \\ - \alpha_j(\bar{u}(\sigma_b) - \bar{u}(\mathbf{y}_{\text{end}}^j)) = m(\sigma_{\text{beg}}^j)g_{\sigma_{\text{beg}}^j} + m(\sigma_{\text{beg}}^j)R_{K, \sigma_{\text{beg}}^j} \end{aligned} \quad (175)$$

(note that, for the sake of simplicity of the notations in (175), we denoted  $\sigma_{\text{beg}}^j = (\sigma_b, \sigma_e)$ ) and on  $\Gamma_j \in \mathcal{S}^-$  and for  $\sigma \in \mathcal{E}_{\text{ext}}$  such that  $(\mathbf{y}_\sigma, \mathbf{y}_{\sigma^+}) \in (\Gamma_j)^2$ , we have



$$\tau_\sigma(\bar{u}(\mathbf{y}_\sigma) - \bar{u}(\mathbf{x}_K)) + \alpha_j(\bar{u}(\mathbf{y}_{\sigma+}) - \bar{u}(\mathbf{y}_\sigma)) + r_{\sigma+}^{j,-} - r_{\sigma}^{j,-} = m(\sigma)g_\sigma + m(\sigma)R_{K,\sigma}, \quad (176)$$

where

$$r_{\sigma}^{j,-} = \alpha_j(\bar{u}(\sigma_b) - \bar{u}(\mathbf{y}_\sigma)), \quad (177)$$

and (the case  $\mathbf{y}_\sigma \in \Gamma_j$  and  $\mathbf{y}_{\sigma+} \notin \Gamma_j$ )

$$\begin{aligned} \tau_\sigma(\bar{u}(\mathbf{y}_{\text{end}}^j) - \bar{u}(\mathbf{x}_K)) + \alpha_j(\bar{u}(\mathbf{y}_{\text{beg}}^j) - \bar{u}(\mathbf{y}_{\text{end}}^j)) + \alpha_j(\bar{u}(\sigma_e) - \bar{u}(\mathbf{y}_{\text{beg}}^j)) \\ - \alpha_j(\bar{u}(\sigma_b) - \bar{u}(\mathbf{y}_{\text{end}}^j)) = m(\sigma_{\text{end}}^j)g_{\sigma_{\text{end}}^j} + m(\sigma_{\text{end}}^j)R_{K,\sigma_{\text{end}}^j} \end{aligned} \quad (178)$$

(note that, for the sake of simplicity of the notations in (178), we denoted  $\sigma_{\text{beg}}^j = (\sigma_b, \sigma_e)$ ). The expansions  $r_{\sigma}^{j,-}$  and  $r_{\sigma}^{j,+}$  can be bounded as

$$|r_{\sigma}^{j,-}| \leq C_{48} \text{size}(\mathcal{T}) \text{ and } |r_{\sigma}^{j,+}| \leq C_{48} \text{size}(\mathcal{T}), \quad (179)$$

where  $C_{48} = \max_{j=1,\dots,N} |\alpha_j| \|\nabla u\|_{L^\infty(\bar{\Omega})}$ ; and in (175) and (178), we have similar estimates, namely (using the second item of Assumption 6.1):

$$|\alpha_j(\bar{u}(\sigma_b) - \bar{u}(\mathbf{y}_{\text{end}}^j))| = |\alpha_j(\bar{u}(\mathcal{S}_j) - \bar{u}(\mathbf{y}_{\text{end}}^j))| \leq C_{48} \text{size}(\mathcal{T}), \text{ where } (\sigma_b, \sigma_e) = \sigma_{\text{beg}}^j, \quad (180)$$

$$|\alpha_j(\bar{u}(\sigma_e) - \bar{u}(\mathbf{y}_{\text{beg}}^j))| = |\alpha_j(\bar{u}(\mathcal{S}_{j-1}) - \bar{u}(\mathbf{y}_{\text{beg}}^j))| \leq C_{48} \text{size}(\mathcal{T}), \text{ where } (\sigma_b, \sigma_e) = \sigma_{\text{end}}^j. \quad (181)$$

(note that  $\mathcal{S}_0 = \mathcal{S}_N$  and  $\{\mathcal{S}_j : j = 1, \dots, N\}$  is the set of the corners of  $\Omega$ ). The equation (63) (in the case of  $\alpha$  constant) remain valid here; this means that

$$- \sum_{\substack{\sigma \in \mathcal{E}_K \\ \sigma = K|L}} \tau_\sigma(e_L - e_K) - \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} \tau_\sigma(e_\sigma - e_K) = - \sum_{\sigma \in \mathcal{E}_K} m(\sigma)R_{K,\sigma}, \forall K \in \mathcal{T}. \quad (182)$$

Subtracting now equation (161) from (173) and (162) from (175), we get, for  $(\mathbf{y}_\sigma, \mathbf{y}_{\sigma-}) \in (\Gamma_j)^2$ , where  $\Gamma_j \in \mathcal{S}^+$

$$\tau_\sigma(e_\sigma - e_K) + \alpha_j(e_\sigma - e_{\sigma-}) + r_{\sigma}^{j,+} - r_{\sigma-}^{j,+} = m(\sigma)R_{K,\sigma}, \quad (183)$$

and, for the case  $\mathbf{y}_\sigma \in \Gamma_j$  and  $\mathbf{y}_{\sigma-} \notin \Gamma_j$

$$\begin{aligned} \tau_\sigma(e_{\text{beg}}^j - e_K) + \alpha_j(e_{\text{beg}}^j - e_{\text{end}}^j) + \alpha_j(\bar{u}(\sigma_e) - \bar{u}(\mathbf{y}_{\text{beg}}^j)) - \alpha_j(\bar{u}(\sigma_b) - \bar{u}(\mathbf{y}_{\text{end}}^j)) \\ = m(\sigma_{\text{beg}}^j)R_{K,\sigma_{\text{beg}}^j}. \end{aligned} \quad (184)$$

We also have similar equalities for  $\Gamma_j \in \mathcal{S}^-$ ; indeed subtracting equation (163) from (176) and (164) from (178), we get, for  $(\mathbf{y}_\sigma, \mathbf{y}_{\sigma+}) \in (\Gamma_j)^2$

$$\tau_\sigma(e_\sigma - e_K) + \alpha_j(e_{\sigma+} - e_\sigma) + r_{\sigma+}^{j,-} - r_{\sigma}^{j,-} = m(\sigma)R_{K,\sigma}, \quad (185)$$

and, for the case  $\mathbf{y}_\sigma \in \Gamma_j$  and  $\mathbf{y}_{\sigma+} \notin \Gamma_j$

$$\begin{aligned} \tau_\sigma(e_{\text{end}}^j - e_K) + \alpha_j(e_{\text{beg}}^j - e_{\text{end}}^j) + \alpha_j(\bar{u}(\sigma_e) - \bar{u}(\mathbf{y}_{\text{end}}^j)) - \alpha_j(\bar{u}(\sigma_b) - \bar{u}(\mathbf{y}_{\text{beg}}^j)) \\ = m(\sigma_{\text{end}}^j)R_{K,\sigma_{\text{end}}^j}. \end{aligned} \quad (186)$$

Furthermore

$$\int_{\Omega} e_{\mathcal{T}}(\mathbf{x}) d\mathbf{x} = 0. \quad (187)$$

Multiplying both sides of the equation (182) by  $e_K, K \in \mathcal{T}$ , summing over  $K, K \in \mathcal{T}$ , using equalities (183)-(186) and Lemma 6.4, we get (recall that  $|(\cdot, \cdot)|_{1, \mathcal{X}(\mathcal{T})}$  and  $|\cdot|_{(\alpha_j), \mathcal{Z}(\mathcal{T})}$  are defined in Definitions 3.4 and 6.3, respectively)

$$|(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}})|_{1, \mathcal{X}(\mathcal{T})}^2 + |\bar{e}_{\mathcal{T}}|_{(\alpha_j), \mathcal{Z}(\mathcal{T})}^2 = \mathbb{T}_5^{\mathcal{T}} + \mathbb{T}_{11}^{\mathcal{T}}, \quad (188)$$

where  $\mathbb{T}_5^{\mathcal{T}}$  is defined by (67) and

$$\begin{aligned} \mathbb{T}_{11}^{\mathcal{T}} &= \sum_{\Gamma_j \in \mathcal{S}^+} \{ -\alpha_j(\bar{u}(\sigma_e) - \bar{u}(\mathbf{y}_{\text{beg}}^j)) + \alpha_j(\bar{u}(\sigma_b) - \bar{u}(\mathbf{y}_{\text{end}}^j)) \} e_{\text{beg}}^j \\ &+ \sum_{\Gamma_j \in \mathcal{S}^+} \sum_{(\mathbf{y}_{\sigma}, \mathbf{y}_{\sigma^-}) \in (\Gamma_j)^2} (-r_{\sigma^-}^{j,+} + r_{\sigma^-}^{j,+}) e_{\sigma} \\ &+ \sum_{\Gamma_j \in \mathcal{S}^-} \{ -\alpha_j(\bar{u}(\sigma_e) - \bar{u}(\mathbf{y}_{\text{beg}}^j)) + \alpha_j(\bar{u}(\sigma_b) - \bar{u}(\mathbf{y}_{\text{end}}^j)) \} e_{\text{end}}^j \\ &+ \sum_{\Gamma_j \in \mathcal{S}^-} \sum_{(\mathbf{y}_{\sigma}, \mathbf{y}_{\sigma^+}) \in (\Gamma_j)^2} (-r_{\sigma^+}^{j,-} + r_{\sigma^+}^{j,-}) e_{\sigma} \\ &+ \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} m(\sigma) R_{K,\sigma} e_{\sigma}. \end{aligned} \quad (189)$$

Reordering the sum, using the fact that  $\bar{u}$  takes the same value on the corners  $\mathcal{S}_j$  (consequence of the second item of Assumption 6.1), assuming that the mesh  $\mathcal{T}$  satisfies the condition (72) and using (179), (180) and (181), we get

$$|\mathbb{T}_{11}^{\mathcal{T}}| \leq C_{49} \sqrt{\text{size}(\mathcal{T})} \left( |(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}})|_{1, \mathcal{X}(\mathcal{T})}^2 + |\bar{e}_{\mathcal{T}}|_{(\alpha_j), \mathcal{Z}(\mathcal{T})}^2 \right)^{\frac{1}{2}}, \quad (190)$$

where  $C_{49}$  depends on  $(\alpha, \zeta_2, u, \Omega)$ .

This with (188), (69)-(70) and techniques of Section 4.4 yields the following Error Estimate result:

**THEOREM 6.6 (ERROR ESTIMATE WHEN  $\alpha$  IS PIECEWISE CONSTANT)** Assume Assumption 6.1 (which gives the discrete compatibility condition (168)). Let  $\mathcal{T}$  be an admissible mesh in the sense of Definition 3.1 and  $\{(f_K, g_{\sigma}), (K, \sigma) \in \mathcal{T} \times \mathcal{E}_{\text{ext}}\}$  be defined by (9). Then, the system (160)-(165) has a unique solution,  $(u_{\mathcal{T}}, v_{\mathcal{T}}) \in \mathcal{X}(\mathcal{T})$ , in the sense of Definition 6.5. Let  $u$  be the solution given by Assumption 6.1 (we can see that  $u = v_x$ , where  $v$  is the solution of the Dirichlet equation (153), satisfies an equation of the form (151) and Assumption 6.1, see Remark 12). Then, there exist  $(C_{50}, C_{51})$ , only depending on  $(\alpha, \zeta_2, u, \Omega)$ , such that:

$$\left( |(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}})|_{1, \mathcal{X}(\mathcal{T})}^2 + |\bar{e}_{\mathcal{T}}|_{(\alpha_j), \mathcal{Z}(\mathcal{T})}^2 \right)^{\frac{1}{2}} \leq C_{50} \sqrt{\text{size}(\mathcal{T})}, \quad (191)$$

and

$$\|e_{\mathcal{T}}\|_{L^2(\Omega)} \leq C_{51} \sqrt{\text{size}(\mathcal{T})}, \quad (192)$$

where  $(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}}) \in \mathcal{X}(\mathcal{T})$  is defined by  $(e_{\mathcal{T}}, \bar{e}_{\mathcal{T}}) = (e_K, e_{\sigma})$ , on  $(K, \sigma)$ , for all  $(K, \sigma) \in \mathcal{T} \times \mathcal{E}_{\text{ext}}$ .

The Error Estimate (191) gives the following approximation for the fluxes  $\int_{\sigma} \nabla u(\mathbf{x}) \cdot \mathbf{n}_{K,\sigma}$ ,  $\sigma \in \mathcal{E}_K$ , for all  $K \in \mathcal{T}$ :

$$\begin{aligned} & \sum_{\sigma=K|L} m(\sigma) d_{\sigma} \left( \frac{u_L - u_K}{d_{\sigma}} - \frac{1}{m(\sigma)} \int_{\sigma} \nabla u(\mathbf{x}) \cdot \mathbf{n}_{K,\sigma} d\gamma(\mathbf{x}) \right)^2 \\ & + \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} m(\sigma) d_{\sigma} \left( \frac{u_{\sigma} - u_K}{d_{\sigma}} - \frac{1}{m(\sigma)} \int_{\sigma} \nabla u(\mathbf{x}) \cdot \mathbf{n}_{K,\sigma} d\gamma(\mathbf{x}) \right)^2 \leq (C_{50})^2 \text{size}(\mathcal{T}). \end{aligned} \quad (193)$$

*Remark 16* (Estimate (192) and Assumption 6.1)

1. The Error Estimate (192) implies the uniqueness of the solution which satisfies (151) and Assumption 6.1.
2. As a consequence of the first items of the Remarks 16 and 13, that, in general, the finite volume solution of the system (160)-(165) does not converge to  $u$  (in the sense of the Error Estimate (192)), when  $u \in \mathcal{C}^2(\overline{\Omega})$  satisfies (151) and only the first item of Assumption 6.1 (without the second item of this Assumption). This includes, for instance, the particular case of the first item of the Remark 13. For this particular case, the finite volume solution  $(u_{\mathcal{T}}, v_{\mathcal{T}}) = (0, 0)$  converges to 0 (which satisfies the Assumption 6.1), in the sense of (192), and not to  $x - y$  (recall that  $x - y$  does not satisfy the second item of the Assumption 6.1).

The authors would like to thank the referees for many useful comments and interesting suggestions. Some corrections on the paper were made at the time of a postdoctoral position of the first author in Weierstrass Institute for Applied Analysis and Stochastics in Berlin.

## References

- [A 06] ATFEH B. AND BRADJI A. , (2006), “Improved convergence order of finite volume solutions. Part II: 2D”. To appear in Arab. J. Math. Sc. .
- [B 05] BRADJI A. , (2005), “Improved convergence order in finite volume and finite element methods”. Thesis, Université de Marseille.
- [BG 05] BRADJI A. AND GALLOUËT T. , (2005), “Finite volume approximation for an oblique derivative boundary problem”. Proceedings of Finite Volume for Complex Applications IV, Hermes. Editors: F. Benkhedhoun, D. Ouazar and S. Raghay, 143-152.
- [EGH 00] EYMARD R., GALLOUËT T. AND HERBIN R. , (2000) “Finite volume methods”. Handbook for Numerical Analysis, Ph. Ciarlet J.L. Lions (Eds.), North Holland, vol. VII pp. 715-1022.
- [GHV 00] GALLOUËT T., HERBIN R. AND VIGNAL M.H. , (2000), “Error estimates for the approximate finite volume solution of convection diffusion equations with general boundary conditions”. SIAM J. Numer. Anal., vol. 37, 1935- 1972.
- [G 85] GRISVARD P. , (1985), “Elliptic problems in non smooth domains”. Pitman Publishing, Monographs and Studies in Mathematics, vol. 24.
- [H 96] HERBIN R. , (1996), “Finite volume methods for diffusion convection equations on general meshes”. Proceedings of Finite Volume for Complex Applications, Problems and Perspectives, Hermes. Editors: F. Benkhedhoun and R. Vilsmeier, 153-160.

- [M 04] MEDKOVÁ D. , (2004), “The oblique derivative problem for the Laplace equation in a plain domain”. *Integral Equations and Operator Theory*, vol. 48, 225-248.
- [M 02] MEHATS F. , (2002), “Convergence of a numerical scheme for a nonlinear oblique derivative boundary value problem”. *Mathematical Modelling and Numerical Analysis*, vol. 36 n° 6, 1111-1132.
- [M 74] MOUSSAOUI M. , (1974), “Régularité de la solution d’un problème à dérivée oblique”. *C.R. Acad. Sci. Paris Sér. A*, 279, 869-872.